Scribed by Luowen Qian

## Lecture 8

*In which we use spectral techniques to find certificates of unsatisfiability for random $k$-SAT formulas.*

## 1  Introduction

Given a random $k$-SAT formula with $m$ clauses and $n$ variables, we want to find a certificate of unsatisfiability of such formula within polynomial time. Here we consider $k$ as fixed, usually equal to 3 or 4. For fixed $n$, the more clauses you have, the more constraints you have, so it becomes easier to show that these constraints are inconsistent. For example, for 3-SAT,

1. In the previous lecture, we have shown that if $m > c_3 \cdot n$ for some large constant $c_3$, almost surely the formula is not satisfiable. But it's conjectured that there is no polynomial time, or even subexponential time algorithms that can find the certificate of unsatisfiability for $m = O(n)$.

2. If $m > c \cdot n^2$ for some other constant $c$, we've shown in the last time that we can find a certificate within polynomial time with high probability that the formula is not satisfiable.

   The algorithm for finding such certificate is shown below.

   ---
   **function** Is-3-SAT-Satisfiable(description of the 3-SAT formula)
       **for** $b_1 \in \{0, 1\}$ **do**
           $x_1 \leftarrow b_1$
           **if** Is-2-SAT-Satisfiable$\left(\text{clauses that contains } x_1 = \overline{b_1}\right)$ **then**
               **return** $\perp$
       **return** UNSATISFIABLE

   ---

   We know that we can solve 2-SATs in linear time, and approximately

   $$\frac{\binom{n-1}{2} \cdot m}{\binom{n}{3} \cdot 2} = \frac{3m}{2n + O(1)} > \frac{3}{2}cn - O(1)$$

clauses contains $x_1 = \overline{b_1}$. Similarly when $c$ is sufficiently large, the 2-SATs will almost surely be unsatisfiable. When a subset of the clauses is not satisfiable, the whole 3-SAT formula is not satisfiable. Therefore we can certify unsatisfiability for 3-SATs with high probability.

In general for $k$-SAT,

1. If $m > c_k \cdot n$ for some large constant $c_k$, almost surely the formula is not satisfiable.

2. If $m > c'_k \cdot n^{k-1}$ for some other constant $c'_k$, we can construct a very similar algorithm.

---
**function** IS-$k$-SAT-SATISFIABLE(description of the $k$-SAT formula)
    **for** $b_1, b_2, ..., b_{k-2} \in \{0,1\}^{k-2}$ **do**
        $x_1, x_2, ..., x_{k-2} \leftarrow b_1, b_2, ..., b_{k-2}$
        **if** Is-2-SAT-Satisfiable(
                clauses that contains $x_1 = \overline{b_1} \vee x_2 = \overline{b_2} \vee ... \vee x_{k-2} = \overline{b_{k-2}}$) **then**
            **return** $\perp$
    **return** UNSATISFIABLE

---

Since for every fixed assignments to the first $k-2$ variables, approximately

$$\frac{\binom{n-k+2}{2}}{\binom{n}{k}2^{k-2}} = \frac{k!}{(n^{k-2} + O(n^{k-3}))2^{k-1}}$$

portion of the $m$ clauses remains, we expect the constant $c'_k = \Omega\left(\frac{2^k}{k!}\right)$ and the running time is $O(2^k m)$.

So what about $m$'s that are in between? It turns out that we can do better with spectral techniques. And the reason that spectral techniques work better is that unlike the previous method, it does not try all the possible assignments and fails to find a certificate of unsatisfiability.

## 2   Reduce certifying unsatisfiability for k-SAT to finding largest independent set

### 2.1   From 3-SAT instances to hypergraphs

Given a random 3-SAT formula $f$, which is an and of $m$ random 3-CNF-SAT clauses over $n$ variables $x_1, x_2, ..., x_n$ (abbreviated as vector $\mathbf{x}$), i.e.

$$f(\mathbf{x}) = \bigwedge_{i=1}^{m} \left(x_{\sigma_{i,1}} = b_{i,1} \vee x_{\sigma_{i,2}} = b_{i,2} \vee x_{\sigma_{i,3}} = b_{i,3}\right),$$

where $\sigma_{i,j} \in [n], b_{i,j} \in \{0,1\}, \forall i \in [m], \sigma_{i,1} < \sigma_{i,2} < \sigma_{i,3}$ and no two $(\sigma_{i,1}, b_{i,1}, \sigma_{i,2}, b_{i,2}, \sigma_{i,3}, b_{i,3})$ are exactly the same. Construct hypergraph $H_f = (X, E)$, where

$$X = \{(i,b) | i \in [n], b \in \{0,1\}\}$$

is a set of $2n$ vertices, where each vertex means an assignment to a variable, and

$$E = \{e_j | j \in [m]\}, e_j = \{(\sigma_{j,1}, \overline{b_{j,1}}), (\sigma_{j,2}, \overline{b_{j,2}}), (\sigma_{j,3}, \overline{b_{j,3}})\}$$

is a set of $m$ 3-hyperedges. The reason we're putting in the negation of $b$ is that a 3-CNF clause evaluates to false if and only if all three subclauses evaluate to false. This will be useful shortly after.

First let's generalize the notion of independent set for hypergraphs.

**Definition 1** *An independent set for hypergraph $H = (X, E)$ is a set $S \subseteq X$ that satisfies $\forall e \in E, e \not\subseteq S$.*

**Proposition 1** *If $f$ is satisfiable, $H_f$ has an independent set of size at least $n$. Equivalently if the largest independent set of $H_f$ has size less than $n$, $f$ is unsatisfiable.*

PROOF: Assume $f$ is satisfiable, let $\mathbf{x} \leftarrow \mathbf{y}$ be a satisfiable assignment, where $\mathbf{y} \in \{0,1\}^n$. Then $S = \{(x_i, y_i) | i \in [n]\}$ is an independent set of size $n$. If not, it means some hyperedge $e_j \subseteq S$, so $\sigma_{j,1} = \overline{b_{j,1}} \land \sigma_{j,2} = \overline{b_{j,2}} \land \sigma_{j,3} = \overline{b_{j,3}}$ and the $j$-th clause in $f$ evaluates to false. Therefore $f$ evaluates to false, which contradicts the fact that $\mathbf{y}$ is a satisfiable assignment. $\square$

We know that if we pick a random graph that's sufficiently dense, i.e. the average degree $d > \ln n$, by spectral techniques we will have a certifiable upper bound on the size of the largest independent set of $O\left(\frac{n}{\sqrt{d}}\right)$ with high probability. So if a random graph has $\Omega(n \log n)$ random edges, we can prove that there's no large independent set with high probability.

But if we have a random hypergraph with $\Omega(n \log n)$ random hyperedges, we don't have any analog of spectral theories for hypergraphs that allow us to do this kind of certification. And from the fact that the problem of certifying unsatisfiability of random formula of $\Omega(n \log n)$ clauses is considered to be hard, we conjecture that there doesn't exist a spectral theory for hypergraphs able to replicate some of the things we are able to do on graphs.

However, what we can do is possibly with some loss, to reduce the hypergraph to a graph, where we can apply spectral techniques.

## 2.2 From 4-SAT instances to graphs

Now let's look at random 4-SATs. Similarly we will write a random 4-SAT formula $f$ as:

$$f(\mathbf{x}) = \bigwedge_{i=1}^{m} \left( x_{\sigma_{i,1}} = b_{i,1} \lor x_{\sigma_{i,2}} = b_{i,2} \lor x_{\sigma_{i,3}} = b_{i,3} \lor x_{\sigma_{i,4}} = b_{i,4} \right),$$

where $\sigma_{i,j} \in [n], b_{i,j} \in \{0,1\}, \forall i \in [m], \sigma_{i,1} < \sigma_{i,2} < \sigma_{i,3} < \sigma_{i,4}$ and no two $(\sigma_{i,1}, b_{i,1}, ..., \sigma_{i,4}, b_{i,4})$ are exactly the same. Similar to the previous construction, but instead of constructing another hypergraph, we will construct just a graph $G_f = (V, E)$, where

$$V = \{(i_1, b_1, i_2, b_2) | i_1, i_2 \in [n], b_1, b_2 \in \{0,1\}\}$$

is a set of $4n^2$ vertices and

$$E = \{e_j | j \in [m]\}, e_j = \{(\sigma_{j,1}, \overline{b_{j,1}}, \sigma_{j,2}, \overline{b_{j,2}}), (\sigma_{j,3}, \overline{b_{j,3}}, \sigma_{j,4}, \overline{b_{j,4}})\}$$

is a set of $m$ edges.

**Proposition 2** *If $f$ is satisfiable, $G_f$ has an independent set of size at least $n^2$. Equivalently if the largest independent set of $H_f$ has size less than $n^2$, $f$ is unsatisfiable.*

PROOF: The proof is very similar to the previous one. Assume $f$ is satisfiable, let $\mathbf{x} \leftarrow \mathbf{y}$ be a satisfiable assignment, where $\mathbf{y} \in \{0,1\}^n$. Then $S = \{(x_i, y_i, x_j, y_j) | i, j \in [n]\}$ is an independent set of size $n^2$. If not, it means some edge $e_j \subseteq S$, so $\sigma_{j,1} = \overline{b_{j,1}} \wedge \sigma_{j,2} = \overline{b_{j,2}} \wedge \sigma_{j,3} = \overline{b_{j,3}} \wedge \sigma_{j,4} = \overline{b_{j,4}}$ and the $j$-th clause in $f$ evaluates to false. Therefore $f$ evaluates to false, which contradicts the fact that $\mathbf{y}$ is a satisfiable assignment. $\square$

From here, we can observe that $G_f$ is not a random graph because some edges are forbidden, for example when the two vertices of the edge has some element in common. But it's very close to a random graph. In fact, we can apply the same spectral techniques to get a certifiable upper bound on the size of the largest independent set if the average degree $d > \ln n$, i.e. if $m = \Omega(n^2 \log n)$, we can certify unsatisfiability with high probability, by upper bounding the size of the largest independent set in the constructed graph.

We can generalize this results for all even $k$'s. For random $k$-SAT where $k$ is even, if $m > c_k n^{k/2} \log n$, we can certify unsatisfiability with high probability, which is better than the previous method which requires $m = \Omega(n^{k-1})$. The same $n^{k/2}(\log n)^{O(1)}$ is achievable for odd $k$, but the argument is significantly more complicated.

## 2.3 Certifiable upper bound for independent sets in modified random sparse graphs

Despite odd $k$'s, another question is that in this setup, can we do better and get rid of the $\log n$ term? This term is coming from the fact that spectral norm break down when the average degree $d < \ln n$. However it's still true that random graph doesn't have any large independent sets even when the average degree $d$ is constant. It's just that the spectral norm isn't giving us good bounds any more, since the spectral norm is at most $O\left(\sqrt{\max d}\right) = O\left(\sqrt{\frac{\log n}{\log \log n}}\right)$. So is there something tighter than spectral bounds that could help us get rid of the $\log n$ term? Could we fix this by removing all the high degree vertices in the random graph?

This construction is due to Feige-Ofek. Given random graph $G \sim G_{n,p}$, where the average degree $d = np$ is some large constant. Construct $G'$ by taking $G$ and removing all edges

incident on nodes with degree higher than $2\bar{d}$ where $\bar{d}$ is the average degree of $G$. We denote $A$ for the adjacency matrix of $G$ and $A'$ for that of $G'$. And it turns out,

**Theorem 3** *With high probability,* $\left\| A' - \frac{d}{n}J \right\| \leq O\left(\sqrt{d}\right)$.

It turns out to be rather difficult to prove. Previously we saw spectral results on random graphs that uses matrix traces to bound the largest eigenvalue. In this case, it's hard to do so because the contribution to the trace of a closed walk is complicated by the fact that edges have dependencies. The other approach is that given random matrix $M$, we will try to upper bound $\|M\| = \max_x \frac{x^T M x}{\|x\|^2}$. A standard way for this is to that for every solution, count the instances of $M$ in which the fixed solution is good, and argue that the number of the fixed solutions is small, which tells us that there's no good solution. The problem here is that the set of solutions is infinitely large. So Feige-Ofek discretize the set of vectors, and then reduce the bound on the quadratic form of a discretized vector to a sum of several terms, each of which has to be carefully bounded.

We always have

$$\max \mathrm{IndSet}(G) \leq \max \mathrm{IndSet}(G') \leq \frac{n}{d} \left\| A' - \frac{d}{n}J \right\|$$

and so, with high probability, we get an $O\left(\frac{n}{\sqrt{d}}\right)$ polynomial time upper bound certificate to the size of the independent set for a $G_{n,d/n}$ random graph. This removes the extra $\log n$ term from our analysis of certificates of unsatisfiability for random $k$-SAT when $k$ is even.

## 3   SDP relaxation of independent sets in random sparse graphs

In order to show a random graph has no large independent sets, a more principled way is to argue that there is some polynomial time solvable relaxation of the problem whose solution is an upper bound of the problem.

Let $\mathrm{SDPIndSet}(G)$ be the optimum of the following semidefinite programming relaxation of the Independent Set problem, which is due to Lovász:

$$\begin{aligned}
\max \quad & \sum_{i \in V} \langle \mathbf{x}_i, \mathbf{x}_0 \rangle \\
s.t. \quad & \\
& \|\mathbf{x}_0\|^2 = 1 \\
& \langle \mathbf{x}_0, \mathbf{x}_i \rangle = \|\mathbf{x}_i\|^2 \quad \forall i \in V \\
& \langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0 \quad \forall (i,j) \in E
\end{aligned}$$

Since it's the relaxation of the problem of finding the maximum independent set, $\max \mathrm{IndSet}(G) \leq \mathrm{SDPIndSet}(G)$ for any graph $G$. And this relaxation has a nice property.

**Theorem 4** *For every $0 < p < 1$, and for every graph $G$, we have*

$$\text{SDPIndSet}(G) \leq \frac{1}{p} \cdot ||pJ - A||$$

*where $J$ is the all-one matrix and $A$ is the adjacency matrix of $G$.*

PROOF: First we note that $\text{SDPIndSet}(G)$ is at most

$$
\begin{aligned}
\max \quad & \sum_{i \in V} \langle \mathbf{x}_i, \mathbf{x}_0 \rangle \\
\text{s.t.} \quad & \\
& ||\mathbf{x}_0||^2 = 1 \\
& \sum_{i \in V} \langle \mathbf{x}_0, \mathbf{x}_i \rangle = \sum_{i \in V} ||\mathbf{x}_i||^2 \\
& \sum_{(i,j) \in E} \langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0
\end{aligned}
$$

and this is equal to

$$
\begin{aligned}
\max \quad & \frac{\left( \sum_{i \in V} \langle \mathbf{x}_i, \mathbf{x}_0 \rangle \right)^2}{\sum_{i \in V} ||\mathbf{x}_i||^2} \\
\text{s.t.} \quad & \\
& ||\mathbf{x}_0||^2 = 1 \\
& \sum_{i \in V} \langle \mathbf{x}_0, \mathbf{x}_i \rangle = \sum_{i \in V} ||\mathbf{x}_i||^2 \\
& \sum_{(i,j) \in E} \langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0
\end{aligned}
$$

which is at most

$$
\begin{aligned}
\max \quad & \frac{\left\| \sum_{i \in V} \mathbf{x}_i \right\|^2}{\sum_{i \in V} ||\mathbf{x}_i||^2} \\
\text{s.t.} \quad & \\
& ||\mathbf{x}_0||^2 = 1 \\
& \sum_{i \in V} \langle \mathbf{x}_0, \mathbf{x}_i \rangle = \sum_{i \in V} ||\mathbf{x}_i||^2 \\
& \sum_{(i,j) \in E} \langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0
\end{aligned}
$$

because

6

$$\sum_{i \in V} \langle \mathbf{x}_i, \mathbf{x}_0 \rangle = \left\langle \sum_{i \in V} \mathbf{x}_i, \mathbf{x}_0 \right\rangle \leq \left\| \sum_{i \in V} \mathbf{x}_i \right\| \cdot \|\mathbf{x}_0\| = \left\| \sum_{i \in V} \mathbf{x}_i \right\|$$

Finally, the above optimization is equivalent to the following

$$\max \quad \frac{\left\| \sum_{i \in V} \mathbf{x}_i \right\|^2 - \frac{1}{p} \sum_{i,j} A_{i,j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle}{\sum_{i \in V} \|\mathbf{x}_i\|^2}$$

$$s.t.$$

$$\|\mathbf{x}_0\|^2 = 1$$

$$\sum_{i \in V} \langle \mathbf{x}_0, \mathbf{x}_i \rangle = \sum_{i \in V} \|\mathbf{x}_i\|^2$$

$$\sum_{(i,j) \in E} \langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0$$

which is at most the unconstrained problem

$$\max \frac{\left\| \sum_{i \in V} \mathbf{x}_i \right\|^2 - \frac{1}{p} \sum_{i,j} A_{i,j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle}{\sum_{i \in V} \|\mathbf{x}_i\|^2} = \max \frac{\sum_{i,j} \left( J - \frac{1}{p} A \right)_{i,j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle}{\sum_{i \in V} \|\mathbf{x}_i\|^2}$$

$$= \lambda_{\max} \left( J - \frac{1}{p} A \right)$$

$$\leq \frac{1}{p} \|pJ - A\|.$$

□

Recall from the previous section that we constructed $G'$ by removing edges from $G$, which corresponds to removing constraints in our semidefinite programming problem, so $\text{SDPIndSet}(G) \leq \text{SDPIndSet}(G') \leq \left\| J - \frac{1}{p} A' \right\|$, which is by theorem 3 at most $O\left( \frac{n}{\sqrt{d}} \right)$ with high probability.

## 4  SDP relaxation of random k-SAT

From the previous section, we get an idea that we can use semidefinite programming to relax the problem directly and find a certificate of unsatisfiability for the relaxed problem.

Given a random $k$-SAT formula $f$:

$$f(\mathbf{x}) = \bigwedge_{i=1}^{m} \bigvee_{j=1}^{k} x_{\sigma_{i,j}} = b_{i,j}$$

$$= \bigwedge_{i=1}^{m} \overline{\bigvee_{j=1}^{k} x_{\sigma_{i,j}} = b_{i,j}}$$

$$= \bigwedge_{i=1}^{m} \bigwedge_{j=1}^{k} \overline{x_{\sigma_{i,j}} = b_{i,j}}.$$

The satisfiability of $f$ is equivalent of the satisfiability of the following equations:

$$x_i^2 = x_i \qquad \forall i \in [n]$$

$$\sum_{i=1}^{m} \left( 1 - \prod_{j=1}^{k} \left( (-1)^{b_{i,j}} x_{\sigma_{i,j}} + b_{i,j} \right) \right) = m$$

Notice that if we expand the polynomial on the left side, there are some of the monomials having degree higher than 2 which prevents us relaxing these equations to a semidefinite programming problem. In order to resolve this, $\forall A \subseteq \mathbf{x}$ and $|A| \le k/2$ we introduce $x_A = \prod_{i \in A} x_i$. Then we can relax all variables to be vectors, i.e.

$$\|\mathbf{x}_\emptyset\|^2 = 1$$

$$\langle \mathbf{x}_A, \mathbf{x}_B \rangle = \langle \mathbf{x}_C, \mathbf{x}_D \rangle \qquad\qquad \forall A \cup B = C \cup D$$

$$\sum_{i=1}^{m} \left( 1 - \prod_{j=1}^{k} \left( (-1)^{b_{i,j}} \mathbf{x}_{\sigma_{i,j}} + b_{i,j} \right) \right) = m \qquad \text{rewritten as quadratic forms of } \mathbf{x}_A$$

For example, if we have a 4-SAT clause

$$x_3 \vee \overline{x_4} \vee x_7 \vee \overline{x_{10}},$$

we can rewrite it as

$$1 - (1 - \mathbf{x}_3) \cdot \mathbf{x}_4 \cdot (1 - \mathbf{x}_7) \cdot \mathbf{x}_{10} = 1 - \mathbf{x}_4 \mathbf{x}_{10} + \mathbf{x}_3 \mathbf{x}_4 \mathbf{x}_{10} + \mathbf{x}_3 \mathbf{x}_7 \mathbf{x}_{10} - \mathbf{x}_3 \mathbf{x}_4 \mathbf{x}_7 \mathbf{x}_{10}$$

$$= 1 - \mathbf{x}_{\{4\}} \mathbf{x}_{\{10\}} + \mathbf{x}_{\{3,4\}} \mathbf{x}_{\{10\}} + \mathbf{x}_{\{3,7\}} \mathbf{x}_{\{10\}} - \mathbf{x}_{\{3,4\}} \mathbf{x}_{\{7,10\}}.$$

For this relaxation, we have:

1. If $m < c(k,n) n^{k/2}$, the SDP associated with the formula is feasible with high probability, where $c(k,n) = 1/n^{o(1)}$ for every fixed $k$.

2. If $m > c'(k) n^{k/2}$, the SDP associated with the formula is not feasible with high probability, where $c'(k,n)$ is a constant for every fixed even $k$, and $c'(k,n) = \text{poly}(\log n)$ for every fixed odd $k$.