

# Lecture 13

## Last time

- ▷ Proof of Davis-Kahan  $\sin \theta$ .
- ▷ Wedin's theorem
- ▷ Community detection

## Today

- ▷ Community detection continued
- ▷ Nets, coverings, and packings.

## Community Detection Continued

In expectation this matrix has a block structured. Assuming the first community is  $[n/2]$ , we have

$$\mathbb{E} A = \left[ \begin{array}{cc|cc} p & \cdots & p & q & \cdots & q \\ \vdots & & \vdots & \vdots & & \vdots \\ p & \cdots & p & q & \cdots & q \\ \hline q & \cdots & q & p & \cdots & p \\ \vdots & & \vdots & \vdots & & \vdots \\ q & \cdots & q & p & \cdots & p \end{array} \right]$$

Check that this matrix is rank

2 with  $\lambda_1 = \frac{p+q}{2} n$ ,  $\lambda_2 = \frac{p-q}{2} n$ ,

$$u_1 = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad \text{and} \quad u_2 = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ -1 \end{bmatrix} \}^{n/2}$$

**Key Insight:** Thus, if  $A$  is close to  $\mathbb{E} A$  we could use its second eigen vector to identify the communities.

# Spectral Clustering Algorithm

Input: Graph  $G$

Step 1: Compute adjacency matrix  $A$

Step 2: Compute  $u_2(A) \leftarrow$  eigenvalue associated with 2nd largest  $\lambda$ ;

Step 3: Return  $\text{sign}(u_2(A))$ .

We shall prove the following theorem

Theorem: Let  $G \sim G(n, p, q)$  and

$q \wedge (p-q) =: \mu > 0$ . Then, with probability at least  $1 - 4e^{-n}$ , the spectral clustering algorithm identifies the communities of  $G$  with at most  $C/\mu^2$  misclassified vertices.

+

Universal constant

Remark:

- In HW3 you will get rid of the dependency in  $q$  from  $\mu$ .
- This result allows for an expected average degree of  $\sqrt{n}$ . This is highly suboptimal, state of the art results

allow for  $O(\log n)$ . See Abbe, Bandeira, & Hall (2015).

Proof: Applying the Davis-Kahan sine theorem

$$\sin \angle(u_2(\mathbb{E}A), u_2(A)) \stackrel{\text{Unit norm eigenvector associated}}{\leq} \frac{2\|A - \mathbb{E}A\|_{\text{op}}}{\min\{\lambda_2(\mathbb{E}A), \lambda_1(\mathbb{E}A) - \lambda_2(\mathbb{E}A)\}}.$$

Check (HW 3)

$$\leq 2 \frac{\|A - \mathbb{E}A\|_{\text{op}}}{\mu n}.$$

Next we use a fact that will occupy the next few lectures.

Fact (00): We have that

$$P(\|A - \mathbb{E}A\|_{\text{op}} \geq C\sqrt{n}) \leq 4e^{-n}$$

for some universal constant  $C > 0$ .

Thus, assuming we are in the event we get

$$\sin \angle(u_2(\mathbb{E}A), u_2(A)) \leq \frac{C\sqrt{n}}{\mu n} \leq \frac{C}{\mu\sqrt{n}}.$$

Therefore,

$$\min_{S \in \{+1\}} \|u_2(EA) - Su_2(A)\|_2 \lesssim \frac{C}{\sqrt{n}}$$

Now, consider  $v_2 := \sqrt{n} u_2$  we get

$$\#\{i \mid \text{sign}(v_2^{(i)}(A)) \neq v_2^{(i)}(EA)\} \xrightarrow{\pm 1 \text{ entries}}$$

$$\begin{aligned} & \leq \sum_i^n \mathbb{1}\{\text{sign}(v_2^{(i)}(A)) \neq v_2^{(i)}(EA)\} \\ & \leq \sum_i^n (v_2^{(i)}(A) - v_2^{(i)}(EA))^2 \\ & \leq C^2/n \end{aligned}$$

□

## Nets, coverings, and packings

Our next goal is to prove Fact(80). Suppose  $A \in \mathbb{R}^{n \times m}$  random matrix and our goal is to get high probability bounds on

$$\|A\|_{\text{op}} = \max_{v \in S^{m-1}} \|Av\|_2.$$

Before, we derived probability bounds  $\max_{i \in [n]} X_i$  using the union bound. However, this only applied for

finitely many variables. In  $\|A\|_{op}$  we have infinitely many of them. We will develop a powerful method to bound maxima  $\max_{v \in K} X_v$  of random variables  $X_v$  that change continuously wrt to the index  $v$ .

**Key idea:** If we substitute the infinite set  $S^{m-1}$  by a finite set  $N \subseteq S^{m-1}$  such that

$$\max_{v \in S^{m-1}} \|Av\|_2 \approx \max_{v \in N} \|Av\|_2.$$

Then, we can apply the same union bound strategy as before. -

In what follows we learn how to construct such finite sets  $N$ .

**Def ( $\epsilon$ -nets):** Let  $(T, d)$  be a metric space. Consider a set  $K \subseteq T$  and a number  $\epsilon > 0$ . A subset  $N \subseteq K$  is called an  $\epsilon$ -net of  $K$  if

$$\forall x \in K \exists x_0 \in N \text{ s.t. } d(x, x_0) \leq \epsilon.$$

The covering number of  $K$  denoted  $N(K, \epsilon, \delta)$  is the cardinality of the smallest  $\epsilon$ -net of  $K$ .

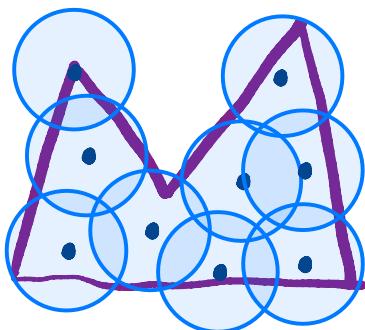
+

Metric space:  $d : T \times T \rightarrow \mathbb{R}_+$  s.t  $\forall x, y, z$

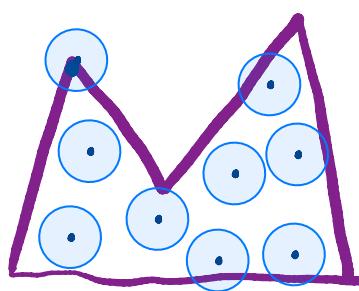
1.  $d(x, x) = 0$ ,
2. If  $x \neq y \Rightarrow d(x, y) > 0$ ,
3.  $d(x, y) = d(y, x)$ ,
4.  $d(x, y) \leq d(x, z) + d(z, y)$ .

Example:  $(\mathbb{R}^d, d)$  with  $d(x, y) = \|x - y\|_2$ .

In that case we cover the set with round balls.



$\epsilon$ -net



$\epsilon/2$ -packing

Def (Packing number): A subset  $N$  of a metric space  $(T, d)$  is  $\epsilon$ -separated if  $d(x, y) > \epsilon$  if  $x, y \in N$  and  $x \neq y$ .

The cardinality of the largest  $\epsilon$ -separa

led subset of  $K \subseteq T$  is called the packing number of  $K$  denoted  $P(K, d, \epsilon)$ .

In turn, covering and packing numbers are almost the same.

**Lemma 0:** For any subset  $K \subseteq T$  and  $\epsilon > 0$ , we have

$$P(K, d, 2\epsilon) \leq N(K, d, \epsilon) \leq P(K, d, \epsilon).$$

**Proof:** To prove the lower bound. Take a  $2\epsilon$ -packing  $P$  and an  $\epsilon$ -covering  $N$ . Let  $x \in P$ , then by def. there is  $y \in N$  s.t.  $d(x, y) \leq \epsilon$ . By def  $\forall z \in P \setminus \{x\}$

$$d(z, y) \geq d(z, x) - d(x, y) > \epsilon.$$

Thus, for each  $x \in P$  there exist a unique  $y \in N$  and so  $|P| \leq |N|$ . Since  $P$  and  $N$  are arbitrary, the lower bound follows.

To prove the upper bound. Let  $N$

be a maximal  $\epsilon$ -separated set of  $K$ , i.e.,  $|N| = P(K, d, \epsilon)$ . We want to show that  $N$  is an  $\epsilon$ -net. Let  $x \in K$ , suppose in search of contradiction that  $\forall y \in N$

$$d(x, y) > \epsilon.$$

This, implies that  $N \cup \{x\}$  is a larger  $\epsilon$ -separated set.  $\varnothing$

Thus,

$$N(K, d, \epsilon) \leq |N| = P(K, d, \epsilon).$$

□