Péter Csikvári

# Applied Discrete Mathematics

Lecture note

# Contents

**Bibliography**                                                                    **71**

# 1.  Spectral Graph Theory

All graphs in this chapter are simple if otherwise not stated, i. e., there are no loops and multiple edges. (Actually, it will be otherwise stated in some proof in the section on the Laplacian-matrix, where it will be convenient to allow multiple edges.)

Apart from the last section on Laplacian-matrix we are concerning with the adjacency matrix of a graph $G$. The adjacency matrix $A(G)$ of a simple graph $G = (V, E)$ is defined as follows: it is a symmetric matrix of size $|V| \times |V|$ labelled by the vertices of the graph $G$, and

$$A(G)_{u,v} = \begin{cases} 1 & \text{if } (u, v) \in E(G), \\ 0 & \text{if } (u, v) \notin E(G). \end{cases}$$

If the graph $G$ is clear from the context we will simple write $A$ instead of $A(G)$.

It is important to understand what it means that we multiply a vector $\underline{x} \in \mathbb{R}^V$ with $A(G)$:

$$(A\underline{x})_u = \sum_{v \in V(G)} A_{u,v} x_v = \sum_{v \in N_G(u)} x_v,$$

where $N_G(u)$ is the set of neighbors of $u$ in the graph $G$. So multiplication with $A(G)$ simply means that we add up the values of the vector on the neighbors of a vertex $u$ and we write this sum in place of $u$. In particular, if $\underline{x}$ is an eigenvector of $A$, i. e., $A\underline{x} = \lambda\underline{x}$ then for all vertex $u$ we have

$$\lambda x_u = \sum_{v \in N_G(u)} x_v.$$

Recall from linear algebra that since $A$ is a real symmetric matrix all eigenvalues are real, and we can choose a basis consisting of orthonormal eigenvectors. Note that if $A\underline{x} = \lambda\underline{x}$ and $A\underline{y} = \mu\underline{y}$ and $\lambda \neq \mu$ then $\underline{x}$ and $\underline{y}$ is immediately orthogonal to each other. If $\lambda = \mu$ it is not necessarily true, but we can still choose orthogonal eigenvectors from this eigenspace.

The goal of this chapter to give a very brief account to spectral graph theory. Here we try to understand the connection between the eigenvalues of the adjacency matrix and the properties of the graph. For instance, it will turn out that the largest eigenvalue is a degree-like concept that is sandwiched between the average degree and the largest degree of the graph. The other eigenvalues of the graph grasp the pseudorandomness and expansion properties of the graph.

Suggested reading:

A. E. Brouwer and W. H. Haemers: Spectra of Graphs [3]. An earlier version of the book is available online.

J. Matousek: *Thirty-three miniatures: mathematical and algorithmic applications of linear algebra* [8]

R. Stanley: *Topics in algebraic combinatorics* [10]

These books are also available online. In the latter case, the online version doesn't contain exercises. Besides the books I would like to call attention to the extensive resources provided by the SageMath software. For more details, see Section 1.5.

## 1.1   Just linear algebra

Many of the things described in this section is just the Frobenius–Perron theory specialized for our case. On the other hand, we will cheat. Our cheating is based on the fact that we will only work with symmetric matrices, and so we can do some shortcuts in the arguments.

We will use the fact many times that if $A$ is a $n \times n$ real symmetric matrix then there exists a basis of $\mathbb{R}^n$ consisting of eigenvectors which we can choose[1] to be orthonormal. Let $\underline{u}_1, \ldots, \underline{u}_n$ be the orthonormal eigenvectors belonging to $\lambda_1 \geq \cdots \geq \lambda_n$: we have $A\underline{u}_i = \lambda_i \underline{u}_i$, and $(\underline{u}_i, \underline{u}_j) = \delta_{ij}$.

Let us start with some elementary observations.

**Proposition 1.1.1.** *If $G$ is a simple graph then the eigenvalues of its adjacency matrix $A$ satisfies $\sum_i \lambda_i = 0$ and $\sum \lambda_i^2 = 2e(G)$, where $e(G)$ denotes the number of edges of $G$. In general, $\sum \lambda_i^\ell$ counts the number of closed walks of length $\ell$.*

---

[1]For the matrix I, any basis will consists of eigenvectors as every vectors are eigenvectors, but of course they won't be orthonormal immediately.

*Proof.* Since $G$ has no loop we have

$$\sum_i \lambda_i = \mathrm{Tr}A = 0.$$

Since $G$ has no multiple edges, the diagonal of $A^2$ consists of the degrees of $G$. Hence

$$\sum_i \lambda_i^2 = \mathrm{Tr}A^2 = \sum d_i = 2e(G).$$

The third statement also follows from the fact $TrA^\ell$ is nothing else than the number of closed walks of length $\ell$. $\square$

Using the next well-known statement, Proposition 1.1.2, one can refine the previous statement such a way that the number of walks of length $\ell$ between vertex $i$ and $j$ can be obtained as

$$\sum_k c_k(i,j)\lambda_k^\ell.$$

The constant $c_k(i,j) = u_{ik}u_{jk}$ if $\underline{u}_k = (u_{1k}, u_{2k}, \dots, u_{nk})$.

**Proposition 1.1.2.** *Let $U = (\underline{u}_1, \dots, \underline{u}_n)$ and $S = \mathrm{diag}(\lambda_1, \dots, \lambda_n)$, then*

$$A = USU^T$$

*or equivalently*

$$A = \sum_{i=1}^n \lambda_i \underline{u}_i \underline{u}_i^T.$$

*Consequently, we have*

$$A^\ell = \sum_{i=1}^n \lambda_i^\ell \underline{u}_i \underline{u}_i^T.$$

*Proof.* First of all, note that $U^T = U^{-1}$ as the vectors $u_i$ are orthonormal. Let $B = USU^T$. Let $\underline{e}_i = (0, \dots, 0, 1, 0, \dots, 0)$, where the i'th coordinate is 1. Then

$$B\underline{u}_i = USU^T u_i = US\underline{e}_i = (\lambda_1\underline{u}_1, \dots, \lambda_n\underline{u}_n)\underline{e}_i = \lambda_i\underline{u}_i = A\underline{u}_i.$$

So $A$ and $B$ coincides on a basis, hence $A = B$. $\square$

Let us turn to the study of the largest eigenvalue and its eigenvector.

**Proposition 1.1.3.** *We have*

$$\lambda_1 = \max_{||\underline{x}||=1} \underline{x}^T A\underline{x} = \max_{\underline{x}\neq 0} \frac{\underline{x}^T A\underline{x}}{||x||^2}.$$

*Further, if for some vector $\underline{x}$ we have $\underline{x}^T A\underline{x} = \lambda_1||\underline{x}||^2$, then $A\underline{x} = \lambda_1\underline{x}$.*

*Proof.* Let us write $\underline{x}$ in the basis $\underline{u}_1, \ldots \underline{u}_n$:

$$\underline{x} = \alpha_1 \underline{u}_1 + \cdots + \alpha_n \underline{u}_n.$$

Then

$$||\underline{x}||^2 = \sum_{i=1}^{n} \alpha_i^2.$$

and

$$\underline{x}^T A \underline{x} = \sum_{i=1}^{n} \lambda_i \alpha_i^2.$$

From this we immediately see that

$$\underline{x}^T A \underline{x} = \sum_{i=1}^{n} \lambda_i \alpha_i^2 \leq \lambda_1 \sum_{i=1}^{n} \alpha_i^2 = \lambda_1 ||\underline{x}||^2.$$

On the other hand,

$$\underline{u}_1^T A \underline{u}_1 = \lambda_1 ||\underline{u}_1||^2.$$

Hence

$$\lambda_1 = \max_{||\underline{x}||=1} \underline{x}^T A \underline{x} = \max_{\underline{x} \neq 0} \frac{\underline{x}^T A \underline{x}}{||\underline{x}||^2}.$$

Now assume that we have $\underline{x}^T A \underline{x} = \lambda_1 ||\underline{x}||^2$ for some vector $\underline{x}$. Assume that $\lambda_1 = \cdots = \lambda_k > \lambda_{k+1} \geq \cdots \geq \lambda_n$, then in the above computation we only have equality if $\alpha_{k+1} = \cdots = \alpha_n = 0$. Hence

$$\underline{x} = \alpha_1 \underline{u}_1 + \cdots + \alpha_k \underline{u}_k,$$

and so

$$A\underline{x} = \lambda_1 \underline{x}.$$

$\square$

**Proposition 1.1.4.** *Let $A$ be a non-negative symmetric matrix. There exists a non-zero vector $\underline{x} = (x_1, \ldots, x_n)$ for which $A\underline{x} = \lambda_1 \underline{x}$ and $x_i \geq 0$ for all $i$.*

*Proof.* Let $\underline{u}_1 = (u_{11}, u_{12}, \ldots, u_{1n})$. Let us consider $\underline{x} = (|u_{11}|, |u_{12}|, \ldots, |u_{1n}|)$. Then $||\underline{x}|| = ||\underline{u}_1|| = 1$ and

$$\underline{x}^T A \underline{x} \geq \underline{u}_1^T A \underline{u}_1 = \lambda_1.$$

Then $\underline{x}^T A \underline{x} = \lambda_1$ and by the previous proposition we have $A\underline{x} = \lambda_1 \underline{x}$. Hence $\underline{x}$ satisfies the conditions. $\square$

**Proposition 1.1.5.** *Let $G$ be a connected graph, and let $A$ be its adjacency matrix. Then*

*(a) If $A\underline{x} = \lambda_1 \underline{x}$ and $\underline{x} \neq \underline{0}$, then no entries of $\underline{x}$ is 0.*

*(b) The multiplicity of $\lambda_1$ is 1.*

*(c) If $A\underline{x} = \lambda_1 \underline{x}$ and $\underline{x} \neq \underline{0}$, then all entries of $\underline{x}$ have the same sign.*

*(d) If $A\underline{x} = \lambda \underline{x}$ for some $\lambda$ and $x_i \geq 0$, where $\underline{x} \neq \underline{0}$, then $\lambda = \lambda_1$.*

*Proof.* (a) Let $\underline{x} = (x_1, \ldots, x_n)$ and $\underline{y} = (|x_1|, \ldots, |x_n|)$. As before we have $||\underline{y}|| = ||\underline{x}||$, and

$$\underline{y}^T A \underline{y} \geq \underline{x}^T A \underline{x} = \lambda_1 ||\underline{x}||^2 = \lambda_1 ||\underline{y}||^2.$$

Hence

$$A\underline{y} = \lambda_1 \underline{y}.$$

Let $H = \{i \mid y_i = 0\}$ and $V \setminus H = \{i \mid y_i > 0\}$. Suppose for contradiction that $H$ is not empty. Note that $V \setminus H$ is not empty either as $\underline{x} \neq \underline{0}$. On the other hand, there cannot be any edge between $H$ and $V \setminus H$: if $i \in H$ and $j \in V \setminus H$ and $(i,j) \in E(G)$, then

$$0 = \lambda_1 y_i = \sum_j a_{ij} y_j \geq y_j > 0$$

contradiction. But if there is no edge between $H$ and $V \setminus H$ then $G$ would be disconnected, which contradicts the condition of the proposition. So $H$ must be empty.

(b) Assume that $A\underline{x}_1 = \lambda_1 \underline{x}_1$ and $A\underline{x}_2 = \lambda_1 \underline{x}_2$, where $\underline{x}_1$ and $\underline{x}_2$ are independent eigenvectors. Note that by part (a), the entries of $\underline{x}_1$ is not 0, so we can choose a constant $c$ such that the first entry of $\underline{x} = \underline{x}_2 - c\underline{x}_1$ is 0. Note that $A\underline{x} = \lambda_1 \underline{x}$ and $\underline{x} \neq 0$ since $\underline{x}_1$ and $\underline{x}_2$ were independent. But then $\underline{x}$ contradicts part (a).

(c) If $A\underline{x} = \lambda_1 \underline{x}$, and $\underline{y} = (|x_1|, \ldots, |x_n|)$ then we have seen before that $A\underline{y} = \lambda_1 \underline{y}$. By part (b), we know that $\underline{x}$ and $\underline{y}$ must be linearly dependent so $\underline{x} = \underline{y}$ or $\underline{x} = -\underline{y}$. Together with part (a), namely that there is no 0 entry, this proves our claim.

(d) Let $A\underline{u}_1 = \lambda_1 \underline{u}_1$. By part (c), all entries have the same sign, we can choose it to be positive by replacing $\underline{u}_1$ with $-\underline{u}_1$ if necessary. Suppose for contradiction that $\lambda \neq \lambda_1$. Note that if $\lambda \neq \lambda_1$ then $\underline{x}$ and $\underline{u}_1$ are orthogonal, but this cannot happen

as all entries of both $\underline{x}$ and $\underline{u}_1$ are non-negative, further they are positive for $\underline{u}_1$, and $\underline{x} \neq \underline{0}$. This contradiction proves that $\lambda = \lambda_1$.

$\square$

So part (c) enables us to recognize the largest eigenvalue from its eigenvector: this is the only eigenvector consisting of only positive entries (or actually, entries of the same sign).

**Proposition 1.1.6.** *(a) Let $H$ be a subgraph of $G$. Then $\lambda_1(H) \leq \lambda_1(G)$.*
*(b) Further, if $G$ is connected and $H$ is a proper subgraph, then $\lambda_1(H) < \lambda_1(G)$.*

*Proof.* (a) Let $\underline{x}$ be an eigenvector of length 1 of the adjacency matrix of $H$ such that it has only non-negative entries. Then

$$\lambda_1(H) = \underline{x}^T A(H)\underline{x} \leq \underline{x}^T A(G)\underline{x} \leq \max_{||\underline{z}||=1} \underline{z}^T A(G)\underline{z} = \lambda_1(G).$$

In the above computation, if $H$ has less number of vertices than $G$, then we complete $\underline{x}$ with 0's in the remaining vertices and we denote the obtained vector with $\underline{x}$ too in order to make sense for $\underline{x}^T A(G)\underline{x}$.

(b) Suppose for contradiction that $\lambda_1(H) = \lambda_1(G)$. Then we have equality everywhere in the above computation. In particular $\underline{x}^T A(G)\underline{x} = \lambda_1(G)$. This means that $\underline{x}$ is eigenvector of $A(G)$ too. Since $G$ is connected $\underline{x}$ must be a (or rather "the") vector with only positive entries by part (a) of the above proposition. But then $\underline{x}^T A(H)\underline{x} < \underline{x}^T A(G)\underline{x}$, a contradiction. $\square$

**Proposition 1.1.7.** *(a) We have $|\lambda_n| \leq \lambda_1$.*
*(b) Let $G$ be a connected graph and assume that $-\lambda_n = \lambda_1$. Then $G$ is bipartite.*
*(c) $G$ is a bipartite graph if and only if its spectrum is symmetric to 0.*

*Proof.* (a) Let $\underline{x} = (x_1, \ldots, x_n)$ be a unit eigenvector belonging to $\lambda_n$, and let $\underline{y} = (|x_1|, \ldots, |x_n|)$. Then

$$|\lambda_n| = |\underline{x}^T A\underline{x}| = \left|\sum a_{ij}x_i x_j\right| \leq \sum a_{ij}|x_i||x_j| = \underline{y}^T A\underline{y} \leq \max_{||\underline{z}||=1} \underline{z}^T A\underline{z} = \lambda_1.$$

(Another solution can be given based on the observation that $0 \leq \text{Tr}A^\ell = \sum \lambda_i^\ell$. If $|\lambda_n| > \lambda_1$ then for large enough odd $\ell$ we get that $\sum \lambda_i^\ell < 0$.)

(b) Since $\lambda_1 \geq \cdots \geq \lambda_n$, the condition can only hold if $\lambda_1 \geq 0 \geq \lambda_n$. Again let $\underline{x} = (x_1, \ldots, x_n)$ be a unit eigenvector belonging to $\lambda_n$, and let $\underline{y} = (|x_1|, \ldots, |x_n|)$. Then

$$\lambda_1 = |\lambda_n| = |\underline{x}^T A \underline{x}| = \left| \sum a_{ij} x_i x_j \right| \leq \sum a_{ij} |x_i| |x_j| = \underline{y}^T A \underline{y} \leq \max_{||\underline{z}||=1} \underline{z}^T A \underline{z} = \lambda_1.$$

We need to have equality everywhere. In particular, $\underline{y}$ is the positive eigenvector belonging to $\lambda_1$, and all $a_{ij} x_i x_j$ have the same signs which can be only negative since $\lambda_n \leq 0$. Hence every edge must go between the sets $V^- = \{i \mid x_i < 0\}$ and $V^+ = \{i \mid x_i > 0\}$. This means that $G$ is bipartite.

(c) First of all, if $G$ is a bipartite graph with color classes $A$ and $B$ then the following is a linear bijection between the eigenspace of the eigenvalue $\lambda$ and the eigenspace of the eigenvalue $-\lambda$: if $A\underline{x} = \lambda \underline{x}$ then let $\underline{y}$ be the vector which coincides with $\underline{x}$ on $A$, and $-1$ times $\underline{x}$ on $B$. It is easy to check that this will be an eigenvector belonging to $-\lambda$.

Next assume that the spectrum is symmetric to 0. We prove by induction on the number of vertices that $G$ is bipartite. Since the spectrum of the graph $G$ is the union of the spectrum of the components there must be a component $H$ with smallest eigenvalue $\lambda_n(H) = \lambda_n(G)$. Note that $\lambda_1(G) = |\lambda_n(G)| = |\lambda_n(H)| \leq \lambda_1(H) \leq \lambda_1(G)$ implies that $\lambda_1(H) = -\lambda_n(H)$. Since $H$ is connected we get that $H$ is bipartite and its spectrum is symmetric to 0. Then the spectrum of $G \backslash H$ has to be also symmetric to 0. By induction $G \setminus H$ must be bipartite. Hence $G$ is bipartite. $\qquad \square$

**Proposition 1.1.8.** *Let $\Delta$ be the maximum degree, and let $\overline{d}$ denote the average degree. Then*

$$\max(\sqrt{\Delta}, \overline{d}) \leq \lambda_1 \leq \Delta.$$

*Proof.* Let $\underline{v} = (1, 1, \ldots, 1)$. Then

$$\lambda_1 \geq \frac{\underline{v}^T A \underline{v}}{||\underline{v}||^2} = \frac{2e(G)}{n} = \overline{d}.$$

If the largest degree is $\Delta$ then $G$ contains $K_{1,\Delta}$ as a subgraph. Hence

$$\lambda_1(G) \geq \lambda_1(K_{1,\Delta}) = \sqrt{\Delta}.$$

Finally, let $\underline{x}$ be an eigenvector belonging to $\lambda_1$. Let $x_i$ be the entry with largest absolute value. Then

$$|\lambda_1| |x_i| = \left| \sum_j a_{ij} x_j \right| \leq \sum_j a_{ij} |x_j| \leq \sum_j a_{ij} |x_i| \leq \Delta |x_i|.$$

7

Hence $\lambda_1 \leq \Delta$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Proposition 1.1.9.** *Let $G$ be a d-regular graph. Then $\lambda_1 = d$ and its multiplicity is the number of components. Every eigenvector belonging to $d$ is constant on each component.*

*Proof.* The first statement already follows from the previous propositions, but it also follows from the second statement so let us prove this statement. Let $\underline{x}$ be an eigenvector belonging to $d$. We show that it is constant on a connected component. Let $H$ be a connected component, and let $c = \max_{i \in V(H)} x_i$, let $V_c = \{i \in V(H) \mid x_i = c\}$ and $V(H) \setminus V_c = \{i \in V(H) \mid x_i < c\}$. If $V(H) \setminus V_c$ were not empty then there exists an edge $(i, j) \in E(H)$ such that $i \in V_c$, $j \in V(H) \setminus V_c$. Then

$$dc = dx_i = \sum_{k \in N(i)} x_k \leq x_j + \sum_{k \in N(i), k \neq j} x_k < c + (d-1)c = dc,$$

contradiction. So $\underline{x}$ is constant on each component. $\qquad\qquad\qquad\qquad\square$

## 1.2 Expanders and pseudorandom graphs

In this section $G$ always will be a $d$–regular graph. The goal of this section is to show how $\lambda_2$ and $\lambda_n$ measures the "randomness" of the graph.

Let $S, T \subseteq V(G)$. Let

$$e(S, T) = |\{(u, v) \in E(G) \mid u \in S, \ v \in T\}|.$$

Note that in the above definition, $S$ and $T$ are not necessarily disjoint. For instance, if $S = T$, then maybe a bit counter intuitively, $e(S, S)$ counts 2 times the number of edges induced by the set $S$. If $G$ were random then we would expect $e(S, T) \approx d\frac{|S||T|}{n}$.

**Theorem 1.2.1.** *Let $G$ be a $d$–regular graph on $n$ vertices with eigenvalues $d = \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$. Let $S, T \subseteq V(G)$ such that $S \cup T = V(G)$ and $S \cap T = \emptyset$. Then*

$$(d - \lambda_2)\frac{|S||T|}{n} \leq e(S, T) \leq (d - \lambda_n)\frac{|S||T|}{n}.$$

Before we start proving this theorem, we need a lemma.

**Lemma 1.2.2.** *Let $A$ be a symmetric matrix with eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$ and corresponding orthornormal eigenvectors $\underline{u}_1, \ldots, \underline{u}_n$. Then*

(a)
$$\min_{\underline{x} \neq 0} \frac{\underline{x}^T A \underline{x}}{||x||^2} = \lambda_n.$$

(b)
$$\max_{\underline{x} \perp \underline{u}_1} \frac{\underline{x}^T A \underline{x}}{||x||^2} = \lambda_2.$$

*Proof.* (a) Let $\underline{x} = \alpha_1 \underline{u}_1 + \cdots + \alpha_n \underline{u}_n$. Then

$$\underline{x}^T A \underline{x} = \sum_{i=1}^{n} \lambda_i \alpha_i^2 \geq \lambda_n \sum_{i=1}^{n} \alpha_i^2 = \lambda_n ||x||^2.$$

On the other hand, $\underline{u}_n^T A \underline{u}_n = \lambda_n ||\underline{u}_n||^2$. This proves part (a).

(b) Again let $\underline{x} = \alpha_1 \underline{u}_1 + \cdots + \alpha_n \underline{u}_n$. Since $\underline{x} \perp \underline{u}_1$ we have $\alpha_1 = (\underline{x}, \underline{u}_1) = 0$. Then

$$\underline{x}^T A \underline{x} = \sum_{i=1}^{n} \lambda_i \alpha_i^2 = \sum_{i=2}^{n} \lambda_i \alpha_i^2 \leq \lambda_2 \sum_{i=1}^{n} \alpha_i^2 = \lambda_2 ||x||^2.$$

On the other hand, $\underline{u}_2^T A \underline{u}_2 = \lambda_2 ||\underline{u}_2||^2$. This proves part (b). $\square$

*Proof of the theorem.* Let $|S| = s$ and $|T| = t$. Let us consider the vector $\underline{x}$ which takes the value $t$ on the vertices of $S$ and the value $-s$ on the vertices of $T$. Then $\underline{x}$ is perpendicular to the all $\underline{1}$ vector, indeed $|S|t - |T|s = 0$. Note that $\underline{u}_1 = \frac{1}{\sqrt{n}}\underline{1}$ so $\underline{x}$ is perpendicular to $\underline{u}_1$. Let us consider

$$\sum_{(i,j) \in E(G)} (x_i - x_j)^2 = d \sum_{i=1}^{n} x_i^2 - 2 \sum_{(i,j) \in E(G)} x_i x_j = d||x||^2 - \underline{x}^T A \underline{x}.$$

First of all, by the lemma we have

$$(d - \lambda_2)||\underline{x}||^2 \leq d||x||^2 - \underline{x}^T A \underline{x} \leq (d - \lambda_n)||\underline{x}||^2.$$

On the other hand,

$$\sum_{(i,j) \in E(G)} (x_i - x_j)^2 = e(S,T)(t - (-s))^2 = e(S,T)(s+t)^2 = e(S,T)n^2.$$

Note that
$$||\underline{x}||^2 = ts^2 + st^2 = st(s+t) = stn.$$

Hence
$$(d - \lambda_2)nst \leq e(S,T)n^2 \leq (d - \lambda_n)nst.$$

9

In other words,
$$(d - \lambda_2)\frac{st}{n} \leq e(S,T) \leq (d - \lambda_n)\frac{st}{n}.$$

□

**Definition 1.2.3.** Let $S \subseteq V(G)$. The set of *neighbors* of $S$ is
$$N(S) = \{u \in V(G) \setminus S \mid \exists v \in S : (u,v) \in E(G)\}.$$

**Definition 1.2.4.** A graph $G$ is called $(n, d, c)$-expander if $|V(G)| = n$, it is $d$–regular and
$$|N(S)| \geq c|S|$$
for every set $S$ satisfying $|S| \leq n/2$.

Intuitively, the larger the $c$, the better your network (your graph $G$) is: if you have a gossip then it spreads in a fast way in a good expander.

**Theorem 1.2.5.** *A $d$–regular graph $G$ on $n$ vertices is an $(n, d, c)$–expander with $c = \frac{d - \lambda_2}{2d}$.*

*Proof.* Let $S \subseteq V(G)$ with $|S| \leq n/2$. Let $T = V(G) \setminus S$, note that $|T| \geq n/2$. Then
$$e(S,T) = e(S, N(S)) \leq d|N(S)|.$$

By Theorem 1.2.1 we have
$$e(S,T) \geq (d - \lambda_2)\frac{|S||T|}{n} \geq (d - \lambda_2)|S|\frac{1}{2}.$$

Hence
$$d|N(S)| \geq \frac{d - \lambda_2}{2}|S|.$$

In other words,
$$|N(S)| \geq \frac{d - \lambda_2}{2d}|S| = c|S|.$$

□

The quantity $d - \lambda_2$ is called *spectral gap*.

Let's see another corollary of Theorem 1.2.1. First we start with a definition.

**Definition 1.2.6.** A set $S \subseteq V(G)$ is called an *independent set* if it induces the empty graph. (In other words, $e(S, S) = 0$.) The size of the largest independent set is denoted by $\alpha(G)$.

**Theorem 1.2.7.** *(Hoffman-Delsarte bound) Let $G$ be a $d$–regular graph on $n$ vertices with eigenvalues $d = \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$. Then*

$$\alpha(G) \leq \frac{-\lambda_n n}{d - \lambda_n}.$$

*Proof.* Let $S$ be the largest independent set, and $T = V(G) \setminus S$. Then $|S| = \alpha(G)$, and $e(S,T) = d|S| = d\alpha(G)$. By Theorem 1.2.1 we have

$$e(S,T) \leq (d - \lambda_n)\frac{|S||T|}{n}.$$

Hence

$$d\alpha(G) \leq (d - \lambda_n)\frac{\alpha(G)(n - \alpha(G))}{n}.$$

By dividing by $\alpha(G)$ and multiplying by $n/(d - \lambda_n)$ we get that

$$\frac{nd}{d - \lambda_n} \leq n - \alpha(G).$$

In other words,

$$\alpha(G) \leq \frac{-\lambda_n n}{d - \lambda_n}.$$

$\square$

The Hoffman-Delsarte bound is surprisingly good in a number of cases. Let's see a bit strange application. A family $\mathcal{F} = \{A_1, A_2, \ldots, A_m\}$ is called *intersecting* if $A_i \cap A_j \neq \emptyset$. Assume that $A_i \subseteq \{1, 2, \ldots, n\}$, and $|A_i| = k$ for all $i$. The question is the following: what's the largest possible intersecting family of $k$-element subsets of $\{1, 2, \ldots, n\}$? If $k > n/2$ then any two $k$-subset is intersecting so the question is trivial. So let us assume that $k \leq n/2$. A good candidate for a large intersecting family is the family $\mathcal{F}_1$ of those subsets which contains the element 1 (or actually any fixed element). Then $|\mathcal{F}_1| = \binom{n-1}{k-1}$. Erdős, Ko and Rado proved that this is indeed the largest possible size of an intersecting family of $k$-element subsets of $\{1, 2, \ldots, n\}$. Actually, they also proved that if $n > 2k$ then an intersecting family of size $\binom{n-1}{k-1}$ must contain a fixed element. For $n = 2k$ this is not true: any family will work where you don't choose a set and its complement at the same time. Now we will only prove the weaker statement that $\binom{n-1}{k-1}$ is an upper bound (and actually we will cheat a bit as we cite a very non-trivial statement).

Let us define the following graph $G$: its vertex set consists of the $k$-element subsets of $\{1, 2, \ldots, n\}$ and two sets are joined by an edge if they are disjoint. This

graph is called the Kneser$(n, k)$ graph. An independent set in this graph is exactly an intersecting family. The following theorem about its spectrum is non-trivial, its proof can be found in C. Godsil and G. Royle: Algebraic graph theory, page 200.

**Theorem 1.2.8.** *The eigenvalues of the Kneser$(n, k)$ graph are*

$$(-1)^i \binom{n - k - i}{k - i},$$

*where $i = 0, \ldots, k$. The multiplicity of $\binom{n-k}{k}$ is 1, otherwise the multiplicity of $(-1)^i \binom{n-k-i}{k-i}$ is $\binom{n}{i} - \binom{n}{i-1}$ if $i \geq 1$.*

Note that the Kneser-graph is $\binom{n-k}{k}$–regular, and according to the previous theorem, its smallest eigenvalue is $-\binom{n-k-1}{k-1}$. Then by the Hoffman-Delsarte bound we have

$$\alpha(\text{Kneser}(n, k)) \leq \frac{\binom{n-k-1}{k-1}\binom{n}{k}}{\binom{n-k}{k} + \binom{n-k-1}{k-1}}.$$

Note that

$$\binom{n - k}{k} = \frac{n - k}{k}\binom{n - k - 1}{k - 1},$$

and so the denominator is

$$\left(\frac{n-k}{k} + 1\right)\binom{n - k - 1}{k - 1} = \frac{n}{k}\binom{n - k - 1}{k - 1}.$$

Hence

$$\frac{\binom{n-k-1}{k-1}\binom{n}{k}}{\binom{n-k}{k} + \binom{n-k-1}{k-1}} = \frac{k}{n}\binom{n}{k} = \binom{n - 1}{k - 1}.$$

Voilá! Honestly this is probably the most complicated proof of the Erdős-Ko-Rado theorem, but there are some similar theorems where the only known proof goes through the eigenvalues of some similarly defined graph.

Now let us turn back to estimating $e(S, T)$ for $d$–regular graphs. The following theorem is called the *expander mixing lemma*.

**Theorem 1.2.9** (Expander mixing lemma). *Let $G$ be a $d$–regular graph on $n$ vertices with eigenvalues $d = \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$. Let $\lambda = \max(|\lambda_2|, \ldots, |\lambda_n|) = \max(|\lambda_2|, |\lambda_n|)$. Let $S, T \subseteq V(G)$, then*

$$\left| e(S, T) - d\frac{|S||T|}{n} \right| \leq \lambda\sqrt{|S||T|}.$$

*Proof.* Let $\chi_S$ and $\chi_T$ be the characteristic vectors of the sets $S$ and $T$: so $\chi_S(u) = 1$ if $u \in S$ and 0 otherwise. Observe that

$$e(S, T) = \chi_S^T A \chi_T.$$

Let us write up $\chi_S$ and $\chi_T$ in the orthonormal basis $\underline{u}_1, \ldots, \underline{u}_n$ of eigenvectors. Note that we can choose $\underline{u}_1$ to be $\frac{1}{\sqrt{n}}\underline{1}$. Let

$$\chi_S = \sum_{i=1}^n \alpha_i \underline{u}_i$$

and

$$\chi_T = \sum_{i=1}^n \beta_i \underline{u}_i.$$

Then

$$\chi_S^T A \chi_T = \sum_{i=1}^n \lambda_i \alpha_i \beta_i.$$

Note that $\alpha_1 = (\chi_S, \underline{u}_1) = \frac{|S|}{\sqrt{n}}$, and similarly $\beta_1 = (\chi_T, \underline{u}_1) = \frac{|T|}{\sqrt{n}}$. Hence

$$\lambda_1 \alpha_1 \beta_1 = d\frac{|S|}{\sqrt{n}}\frac{|T|}{\sqrt{n}} = d\frac{|S||T|}{n}.$$

Hence

$$e(S, T) - d\frac{|S||T|}{n} = \sum_{i=2}^n \lambda_i \alpha_i \beta_i.$$

Then

$$\left| e(S, T) - d\frac{|S||T|}{n} \right| = \left| \sum_{i=2}^n \lambda_i \alpha_i \beta_i \right| \le \lambda \sum_{i=2}^n |\alpha_i||\beta_i|.$$

Now let us apply a Cauchy-Schwartz inequality:

$$\sum_{i=2}^n |\alpha_i||\beta_i| \le \left( \sum_{i=2}^n |\alpha_i|^2 \right)^{1/2} \left( \sum_{i=2}^n |\beta_i|^2 \right)^{1/2}.$$

We will be a bit generous:

$$\left( \sum_{i=2}^n |\alpha_i|^2 \right)^{1/2} \left( \sum_{i=2}^n |\beta_i|^2 \right)^{1/2} \le \left( \sum_{i=1}^n |\alpha_i|^2 \right)^{1/2} \left( \sum_{i=1}^n |\beta_i|^2 \right)^{1/2} =$$

$$= ||\chi_S|| \cdot ||\chi_T|| = |S|^{1/2}|T|^{1/2}.$$

Hence

$$\left| e(S, T) - d\frac{|S||T|}{n} \right| \le \lambda\sqrt{|S||T|}.$$

$\square$

If we were not generous at the last step then we could have proved the following stronger statement:

$$\left| e(S,T) - d\frac{|S||T|}{n} \right| \leq \lambda \left( |S| - \frac{|S|^2}{n} \right)^{1/2} \left( |T| - \frac{|T|^2}{n} \right)^{1/2}.$$

We could have used that $\alpha_1 = \frac{|S|}{\sqrt{n}}$ and $\beta_1 = \frac{|T|}{\sqrt{n}}$.

**Remark 1.2.10.** A graph is called $(n, d, \lambda)$-pseudorandom if it is a $d$–regular graph on $n$ vertices with $\max(|\lambda_2|, |\lambda_n|) \leq \lambda$. (Note that for bipartite $d$–regular graphs it is convenient to require that $\lambda_2 \leq \lambda$, we will not do it though.) Many theorems which assert that "a random $d$–regular graph satisfies property $P$ with very high probability" have an analogue that "an $(n, d, \lambda)$-pseudorandom graph with $\lambda \leq ...$ satisfies property $P$". Such an example is the following theorem due to F. Chung.

**Theorem 1.2.11.** *Let $G$ be an $(n, d, \lambda)$–pseudorandom graph. Then the diameter of $G$ is at most*

$$\left\lceil \frac{\log(n-1)}{\log\left(\frac{d}{\lambda}\right)} \right\rceil + 1.$$

*Proof.* We need to prove that there exists an $r \leq \lceil \frac{\log(n-1)}{\log\left(\frac{d}{\lambda}\right)} \rceil + 1$ such that the distance between any vertices $i$ and $j$ is at most $r$. In other words, there is a walk of length at most $r$ starting at vertex $i$ and ending at vertex $j$. It means that we have to prove that $(A^r)_{ij} > 0$. On the other hand, we know that

$$(A^r)_{ij} = \sum_{k=1}^{n} u_{ik} u_{jk} \lambda_k^r,$$

where $\underline{u}_k = (u_{1k}, \ldots, u_{nk})$. As usual, $\underline{u}_1, \ldots, \underline{u}_n$ is an orthonormal basis of eigenvectors: $A\underline{u}_i = \lambda_i \underline{u}_i$, and $\underline{u}_1 = \frac{1}{\sqrt{n}}\underline{1}$. Then

$$u_{i1} u_{j1} \lambda_1^r = \frac{d^r}{n}.$$

So it is enough to prove that

$$\left| \sum_{k=2}^{n} u_{ik} u_{jk} \lambda_k^r \right| < \frac{d^r}{n}$$

for some $r \leq \lceil \frac{\log(n-1)}{\log\left(\frac{d}{\lambda}\right)} \rceil + 1$.

$$\left| \sum_{k=2}^{n} u_{ik} u_{jk} \lambda_k^r \right| \leq \lambda^r \sum_{k=2}^{n} |u_{ik}||u_{jk}| \leq \lambda^r \left( \sum_{k=2}^{n} |u_{ik}|^2 \right)^{1/2} \left( \sum_{k=2}^{n} |u_{jk}|^2 \right)^{1/2} =$$

14

$$= \lambda^r \left( \sum_{k=1}^{n} |u_{ik}|^2 - u_{i1}^2 \right)^{1/2} \left( \sum_{k=2}^{n} |u_{jk}|^2 - u_{i1}^2 \right)^{1/2} =$$

$$= \lambda^r \left( 1 - \frac{1}{n} \right)^{1/2} \left( 1 - \frac{1}{n} \right)^{1/2} = \lambda^r \left( 1 - \frac{1}{n} \right).$$

The second inequality is a Cauchy-Schwartz. After that we used that the row(!) vectors of $U = (\underline{u}_1, \ldots \underline{u}_n)$ have length 1. This is true since the orthonormality of the column vectors implies the orthonormality of the row vectors. (Indeed, $U \cdot U^T = I$ implies $U^T \cdot U = I$.) Note that

$$\lambda^r \left( 1 - \frac{1}{n} \right) < \frac{d^r}{n}$$

indeed holds true for $r = \lceil \frac{\log(n-1)}{\log\left(\frac{d}{\lambda}\right)} \rceil + 1$. □

It is clear from the previous theorems that the smaller the $\lambda$, the better pseudo-random properties $G$ have. Then the following question naturally arises: what's the best $\lambda$ we can achieve? The complete graph $K_{d+1}$ has eigenvalues $d, (-1)^{(d)}$, but the problem is that it is only one graph. What happens if we require our graph to be large? The Alon-Boppana theorem asserts that in some sense $2\sqrt{d-1}$ is a threshold:

**Theorem 1.2.12** (Alon-Boppana). *Let $(G_n)$ be a sequence of d-regular graphs such that $|V(G_n)| \to \infty$. Then*

$$\liminf_{n \to \infty} \lambda_2(G_n) \geq 2\sqrt{d-1}.$$

*In other words, if $s < 2\sqrt{d-1}$ then there are only finitely many d–regular graphs for which $\lambda_2 \leq s$.*
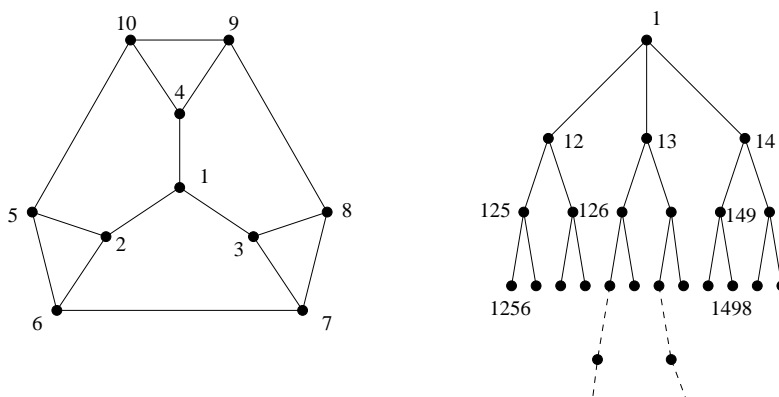
We will prove a slightly stronger statement due to Serre.

**Theorem 1.2.13** (Serre). *For every $\varepsilon > 0$, there exists a $c = c(\varepsilon, d)$ such that for any d–regular graph $G$, the number of eigenvalues $\lambda$ with $\lambda \geq (2 - \varepsilon)\sqrt{d-1}$ is at least $c|V(G)|$.*

Serre's theorem indeed implies the Alon-Boppana theorem since for any $s < 2\sqrt{d-1}$ we choose $\varepsilon$ such that $s < (2 - \varepsilon)\sqrt{d-1}$, then if $|V(G)| > 2/c(\varepsilon, d)$, we have at least two eigenvalues which are bigger then $s$ (one of them is $d$), so $\lambda_2(G) > s$. The following proof of Serre's theorem is due to S. Cioaba.

15

*Proof.* The idea of the proof is that $p_{2k} = \sum_{i=1}^{n} \lambda_i^{2k}$ cannot be too small. Recall that $p_{2k}$ counts the number of closed walks of length $2k$. We will show that for any vertex $v$, the number of closed walks $W_{2k}(v)$ of length $2k$ starting and ending at $v$ is at least as large as the number of closed walks starting and ending at some root vertex of the infinite $d$–regular tree $\mathbb{T}_d$.

Let us consider the following infinite $d$-regular tree, its vertices are labeled by the walks starting at the vertex $v$ which never steps immediately back to a vertex from where it came. Such walks are called non-backtracking walks. For instance, 149831 is such a walk, but 1494 is not a backtracking walk since after 9 we immediately stepped back to 4. We connect two non-backtracking walks in the tree if one of them is a one-step extension of the other.



Note that every closed walk in the tree corresponds to a closed walk in the graph: for instance, $1, 14, 149, 14, 1$ corresponds to $1, 4, 9, 4, 1$. (In some sense, these are the "genuinely" closed walks.) On the other hand, there are closed walks in the graph $G$, like 149831, which are not closed anymore in the tree. Let $r_{2k}$ denote the number of closed closed walks from a given a root vertex in the infinite $d$–regular tree. So far we know that

$$p_{2k} = \sum_{v \in V(G)} W_{2k}(v) \geq n r_{2k}.$$

We would be able to determine $r_{2k}$ explicitly, but for our purposes, it is better to give a lower bound with which we can count easily. Such a lower bound is

$$r_{2k} \geq \frac{\binom{2k}{k}}{k+1} d(d-1)^{k-1} > \frac{1}{(k+1)^2}(2\sqrt{d-1})^{2k}.$$

The second inequality comes from Stirling's formula, so we only need to understand the first inequality. Every closed walk in the tree can be encoded as follows: we write

16

a 1 if we step down (so away from the root) and $-1$ if we step up (towards the root), additionally we choose a direction $d - 1$ or $d$ ways if we step down. More precisely, we can choose our step in $d$ ways if we are in the root and $d - 1$ ways otherwise. (The lower bound $d - 1$ would be sufficient for us.) Note that we have to step down exactly $k$ times, and step up exactly $k$ times to get a closed walk. So the sequence of "directions" is at least $d(d - 1)^{k-1}$. The sequence of $\pm 1$ has two conditions: (i) there must be exactly $k$ 1's and exactly $k$ $-1$'s, (ii) the sum of the first few elements cannot be negative (we cannot go higher than the root): $s_1 + s_2 + \cdots + s_i \geq 0$ for all $1 \leq i \leq 2k$, where $s_i = \pm 1$ according to the i-th step goes down or up. Such sequences are counted by the Catalan-numbers: $\frac{\binom{2k}{k}}{k+1}$.

Now let us finish the proof by using the fact that

$$p_{2k} \geq \frac{n}{(k+1)^2}(2\sqrt{d-1})^{2k}$$

for every $k$. Let $m$ be the number of eigenvalues which are at least $(2 - \varepsilon)\sqrt{d-1}$. Let us consider the sum

$$\sum_{i=1}^{n}(d + \lambda_i)^{2t},$$

where $t$ is a positive integer that we will choose later. Note that $0 \leq d + \lambda_i \leq 2d$, hence

$$\sum_{i=1}^{n}(d + \lambda_i)^{2t} \leq m(2d)^{2t} + (n - m)(d + (2 - \varepsilon)\sqrt{d-1})^{2t}.$$

On the other hand, by the binomial theorem we have

$$\sum_{i=1}^{n}(d + \lambda_i)^{2t} = \sum_{i=1}^{n}\sum_{j=0}^{2t}\binom{2t}{j}d^j\lambda_i^{2t-j} = \sum_{j=0}^{2t}\binom{2t}{j}d^j\left(\sum_{i=1}^{n}\lambda_i^{2t-j}\right).$$

We know that $p_k = \sum_{i=1}^{n}\lambda_i^k \geq 0$ if $k$ is odd and $p_{2k} \geq \frac{n}{(k+1)^2}(2\sqrt{d-1})^{2k}$. Hence

$$\sum_{j=0}^{2t}\binom{2t}{j}d^j\left(\sum_{i=1}^{n}\lambda_i^{2t-j}\right) \geq \sum_{j=0}^{t}\binom{2t}{2j}d^j\left(\sum_{i=1}^{n}\lambda_i^{2t-2j}\right) \geq$$

$$\geq \sum_{j=0}^{t}\binom{2t}{2j}d^j\frac{n}{(t-j+1)^2}(2\sqrt{d-1})^{2t-2j} \geq \frac{n}{(t+1)^2}\sum_{j=0}^{t}\binom{2t}{2j}d^j(2\sqrt{d-1})^{2t-2j} =$$

$$= \frac{n}{2(t+1)^2}\left((d + 2\sqrt{d-1})^{2t} + (d - 2\sqrt{d-1})^{2t}\right) \geq \frac{n}{2(t+1)^2}(d + 2\sqrt{d-1})^{2t}.$$

Hence we have

$$m(2d)^{2t} + (n-m)(d + (2-\varepsilon)\sqrt{d-1})^{2t} \geq \frac{n}{2(t+1)^2}(d + 2\sqrt{d-1})^{2t}.$$

This means that

$$\frac{m}{n} \geq \frac{\frac{1}{2(t+1)^2}(d + 2\sqrt{d-1})^{2t} - (d + (2-\varepsilon)\sqrt{d-1})^{2t}}{(2d)^{2t} - (d + (2-\varepsilon)\sqrt{d-1})^{2t}}.$$

Note that

$$\left(\frac{d + 2\sqrt{d-1}}{d + (2-\varepsilon)\sqrt{d-1}}\right)^{2t}$$

grows much faster than $2(t+1)^2$, so we can choose a $t_0$ for which

$$\frac{1}{2(t_0+1)^2}(d + 2\sqrt{d-1})^{2t_0} - (d + (2-\varepsilon)\sqrt{d-1})^{2t_0} > 0,$$

then

$$c(\varepsilon, d) = \frac{\frac{1}{2(t_0+1)^2}(d + 2\sqrt{d-1})^{2t_0} - (d + (2-\varepsilon)\sqrt{d-1})^{2t_0}}{(2d)^{2t_0} - (d + (2-\varepsilon)\sqrt{d-1})^{2t_0}}$$

satisfies the conditions of the theorem. $\qquad\square$

**Remark 1.2.14.** A $d$–regular non-bipartite graph $G$ is called Ramanujan if $\lambda_2, |\lambda_n| \leq 2\sqrt{d-1}$. If $G$ is bipartite then it is called Ramanujan if $\lambda_2 \leq 2\sqrt{d-1}$. It is known that for any $d$ there exists infinitely many $d$–regular bipartite Ramanujan-graph, this is a result of A. Marcus, D. Spielman and N. Srivastava. On the other hand, if $G$ is non-bipartite then our knowledge is much more limited: for $d = p^\alpha + 1$, where $p$ is a prime there exists construction for infinite family of $d$–regular Ramanujan-graphs. It is conjectured that a random $d$–regular graph is Ramanujan with positive probability independently of the number of vertices.

## 1.3 Derandomization

Suppose we have a boolean function $f : \{0,1\}^n \to \{0,1\}$ and a probabilistic algorithm $\mathbb{A}$ that approximates $f$ using a random $r \in \{0,1\}^m$ in the following sense:

$$\mathbb{P}(\mathbb{A}(x,r) \neq f(x)) \leq \frac{1}{4}.$$

We can significantly reduce the error rate by applying the algorithm $2t + 1$ times and then taking the majority vote. But this requires $(2t + 1)m$ random bits, and

18

unfortunately, random bits have costs. With the help of expanders we can achieve that we only need to generate $m$ random bits, but still reducing the error rate as follows. Take a pseudo-random graph with parameters $(2^m, d, \lambda)$. Pick a random vertex $v \in V$ and output

$$\mathbb{B}(x, v) := \text{Majority}_{u \in N_G(v)} \mathbb{A}(x, u).$$

**Proposition 1.3.1.** *For every $x \in \{0, 1\}^n$ we have*

$$\mathbb{P}(\mathbb{B}(x, v) \neq f(x)) \leq 4 \left( \frac{\lambda}{d} \right)^2.$$

*Proof.* Fix an input $x \in \{0, 1\}^n$. Let

$$S = \{v \in V(G) \mid \mathbb{B}(x, v) \neq f(x)\},$$

and

$$T = \{v \in V(G) \mid \mathbb{A}(x, v) \neq f(x)\}.$$

Then $|T| \leq \frac{1}{4} 2^m$ by the assumption on the algorithm $\mathbb{A}$. Note that every $v \in S$ has at least $d/2$ neighbors in $T$. Hence $e(S, T) \geq \frac{d}{2}|S|$. Expander mixing lemma claims that

$$\left| e(S, T) - d\frac{|S||T|}{2^m} \right| \leq \lambda\sqrt{|S||T|}.$$

Putting these together we get that

$$\lambda\sqrt{|S||T|} \geq e(S, T) - d\frac{|S||T|}{2^m} \geq \frac{d}{2}|S| - d\frac{|S||T|}{2^m} \geq \frac{d}{2}|S| - \frac{d}{4}|S| = \frac{d}{4}|S|.$$

Hence

$$\frac{16\lambda^2}{d^2}|T| \geq |S|.$$

Thus

$$\mathbb{P}(\mathbb{B}(x, v) \neq f(x)) = \frac{|S|}{2^m} \leq \frac{16\lambda^2}{d^2}\frac{|T|}{2^m} \leq 4 \left( \frac{\lambda}{d} \right)^2.$$

$\square$

**Remark 1.3.2.** From the proof we can actually see that

$$\mathbb{P}(\mathbb{B}(x, v) \neq f(x)) \leq 16 \left( \frac{\lambda}{d} \right)^2 \mathbb{P}(\mathbb{A}(x, v) \neq f(x))$$

as long as $\mathbb{P}(\mathbb{A}(x, v) \neq f(x)) \leq 1/4$. It is not always a good idea to reduce the error rate this way since we have to do the deterministic computation of $\mathbb{A}(x, r)$ $d$ times. So there is a trade-off between error rate, running time of $\mathbb{A}$ and the cost of generating a random bit.

## 1.4 Google Page Rank

We have seen that if $A$ is a non-negative (symmetric) matrix, then it has an eigenvector with only non-negative entries. These entries can be used to rank the vertices. Google uses a similar idea to rank websites, but they have a directed graph.

So suppose that $u_1, \ldots, u_m$ are web pages that link to a page $v$. Let page $u_i$ have out-degree $d_i$. Then the Page Rank $w(\cdot)$ satisfies the equation

$$w(v) = 1 - \alpha + \alpha \sum_{i=1}^{m} \frac{w(u_i)}{d_i},$$

where $\alpha$ is a preset constant. Google used $\alpha = 0.85$.

To translate this to the language of linear algebra let $A$ be the adjacency matrix of the directed graph $G$, that is, $A_{uv} = 1$ if there is a directed edge $(u, v) \in E(G)$. Let $D$ be the diagonal matrix containing the out-degrees, and $J$ be the all-one matrix. Set $M = \frac{1-\alpha}{n} J + \alpha D^{-1} A$. (If $d_u = 0$ for some $u$, then we artificially set the value $d_u$ to be 1.) Then $M$ has only positive entries, so it has a unique positive left eigenvector $\underline{w}$ normalized such a way that $\sum_{u \in V(G)} w(u) = 1$. We have $M\underline{1} = \underline{1}$ and $\underline{w}M = \underline{w}$.

The matrix $M$ is huge, but $A$ is sparse, so the sequence $\underline{w}_{k+1} = \underline{w}_k M$ will converge to $\underline{w}$ and the computational cost is not terrible. The value of $\alpha$ regulates the speed of convergence.

## 1.5 SageMath and spectral graph theory

It is always very instructive to see a lot of examples. SageMath, an easy-to-use math software, provides a simple way to compute eigenvalues and eigenvectors of graphs. You don't even need to download this program, you can use the online version at `https://sagecell.sagemath.org/`. It is, of course, possible to download the program.

Let's see an example. Copy-paste the following code into the window and push the Evaluate button under the window:

```
g=graphs.PetersenGraph()
g.show()
print g.adjacency_matrix()
print g.spectrum()
print g.eigenvectors()
```

For the method g.eigenvectors() SageMath gives the eigenvalues again with a basis of the corresponding eigenspace and its dimension. Instead of the Petersen graph, you can try many implemented graphs, see a list at `http://doc.sagemath.org/html/en/reference/graphs/sage/graphs/graph_generators.html` or you can even build up one:

```
g=Graph({})
g.add_edges([(0,1),(1,2),(2,3),(3,4),(2,4)])
g.show()
g.spectrum()
```

Of course, SageMath provides many other tools to work with graphs. For a list, see `http://doc.sagemath.org/html/en/reference/graphs/sage/graphs/generic_graph.html`

# 2.  Probabilistic Method

## 2.1  Introduction

In mathematics we often faces with the problem that we need to prove that a certain structure $S$ exists with a required property $P$. In most cases we simply prove the existence by constructing the required structure $S$. Unfortunately, sometimes this route does not work and we can only give an existence proof, a proof that does not give much besides the existence. A popular tool providing such proofs is for instance the pigeonhole principle. Another tool that we describe in this chapter is the so-called probabilistic method. Using this method we show that in a certain probability space the required structure $S$ exists with positive probability.

In this chapter we consider the most basic methods of this general idea. The first method will simply rely on the union bound. A tiny bit more tricky method is the so-called first moment method. Then we develop further this method: altered first moment method. Finally we will study the so-called second moment method to investigate threshold functions of random graphs.

These methods seem to be very simple, but in reality there are two (non-independent) problems that can occur. First, we have to realize that we need to use the probabilistic method, and forget the idea of constructing the required structure. Second, the construction of the probability space can be very tricky. Below we will see various examples. Some of them will be very simple, others will be quite tricky.

Finally, I would like to call attention to the book The Probabilistic Method by Noga Alon and Joel Spencer [2]. This is a very nicely written book containing all the results of this chapter and many many more beautiful applications of the probabilistic method.

## 2.2 Basics

In this section we collected some notations and basic inequalities.

The expected value of a random variable $X$ is $\int_\Omega X dP$. If $X$ takes only non-negative integers then

$$\mathbb{E}X = \sum_{k=0}^{\infty} k\mathbb{P}(X = k).$$

The variance of a random variable is

$$\mathrm{Var}(X) = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2.$$

The covariance of the random variables $X$ and $Y$ will be denoted by

$$\mathrm{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}X \cdot \mathbb{E}Y.$$

If $X = X_1 + X_2 + \cdots + X_n$ then

$$\mathrm{Var}(X) = \sum_{i=1}^{n} \mathrm{Var}(X_i) + \sum_{i \neq j} \mathrm{Cov}(X_i, X_j).$$

### 2.2.1 Useful inequalities

**Proposition 2.2.1.** *For all $x \in \mathbb{R}$ we have $1 + x \le e^x$.*

*Proof.* If $x > 0$ then

$$1 + x \le \sum_{k=0}^{\infty} \frac{x^k}{k!} = e^x.$$

If $x \le -1$ then the claim is trivial. If $-1 \le x \le 0$, then set $y = -x \ge 0$. Then

$$\frac{1}{1-y} = \sum_{k=0}^{\infty} y^k \ge \sum_{k=0}^{\infty} \frac{y^k}{k!} = e^y.$$

Hence $e^x = e^{-y} \ge 1 - y = 1 + x$. $\qquad\square$

**Proposition 2.2.2.** *We have*

$$\binom{n}{k} \le \left(\frac{en}{k}\right)^k.$$

*Proof.* Since $\binom{n}{k} \le \frac{n^k}{k!}$, it is enough to prove that $k! \ge \left(\frac{k}{e}\right)^k$. This is indeed true:

$$e^k \ge \prod_{j=1}^{k-1}\left(1 + \frac{1}{j}\right)^j = \prod_{j=1}^{k-1}\frac{(j+1)^j}{j^j} = \frac{k^{k-1}}{(k-1)!} = \frac{k^k}{k!}.$$

$\qquad\square$

## 2.2.2 Basic inequalities in probability theory

We recall some basic inequalities.

**Proposition 2.2.3** (Union bound).

$$\mathbb{P}\left(\bigcup_{i=1}^{m} A_i\right) \leq \sum_{i=1}^{m} \mathbb{P}(A_i).$$

**Theorem 2.2.4** (Markov's inequality). *Let $X$ be a non-negative random variable with $\mathbb{E}X > 0$. Then for arbitrary positive $\lambda$ we have*

$$\mathbb{P}(X \geq \lambda) \leq \frac{\mathbb{E}X}{\lambda}.$$

*Proof.*

$$\mathbb{E}X = \int X dP \geq \int_{\{X \geq \lambda\}} X dP \geq \int_{\{X \geq \lambda\}} \lambda dP = \lambda \mathbb{P}(X \geq \lambda).$$

$\square$

A simple corollary of Markov's inequality is Chebyshev's inequality.

**Theorem 2.2.5** (Chebyshev's inequality). *Let $X$ be a random variable with $\mathbb{E}X = \mu$ and $\mathrm{Var}(X) = \sigma^2$. Then*

$$\mathbb{P}(|X - \mu| \geq \lambda\sigma) \leq \frac{1}{\lambda^2}.$$

*Proof.* Let us apply Markov's inequlity to the random variable $Y = (X - \mu)^2$. Then $\mathbb{E}Y = \mathrm{Var}(X) = \sigma^2$ by definition.

$$\mathbb{P}(|X - \mu| \geq \lambda\sigma) = \mathbb{P}(Y \geq \lambda^2\sigma^2) \leq \frac{\mathbb{E}Y}{\lambda^2\sigma^2} = \frac{1}{\lambda^2}.$$

$\square$

$$\star \quad \star \quad \star$$

In this lecture note we would like to prove combinatorial theorems and so many times the studied random variable $X$ takes only non-negative integer values. In fact, often one can translate the combinatorial statement to a statement about the probability that a random variable takes the value 0. This motivates us to collect some results on estimates of $\mathbb{P}(X = 0)$.

**Theorem 2.2.6.** *If $X$ takes only non-negative integer values then*

$$\mathbb{P}(X = 0) \geq 1 - \mathbb{E}X.$$

*Proof.* We have

$$\mathbb{P}(X > 0) = \sum_{k=1}^{\infty} \mathbb{P}(X = k) \leq \sum_{k=0}^{\infty} k\mathbb{P}(X = k) = \mathbb{E}X,$$

or equivalently,

$$\mathbb{P}(X = 0) \geq 1 - \mathbb{E}X.$$

$\square$

This implies that, for instance, if a sequence of random variables $X_n$ satisfies that $\lim_{n\to\infty} \mathbb{E}X_n = 0$, then

$$\lim_{n\to\infty} \mathbb{P}(X_n = 0) = 1.$$

However, $\mathbb{E}X_n \to \infty$ does not guarantee that

$$\lim_{n\to\infty} \mathbb{P}(X_n = 0) = 0.$$

To phrase such a statement we also need the variance of the random variables $X_n$.

**Theorem 2.2.7.**
$$\mathbb{P}(X = 0) \leq \frac{\text{Var}(X)}{(\mathbb{E}X)^2}.$$

*Proof.* Let us use Chebyshev-inequality.

$$\mathbb{P}(X = 0) \leq \mathbb{P}(|X - \mathbb{E}X| \geq \mathbb{E}X) \leq \frac{\text{Var}(X)}{(\mathbb{E}X)^2}.$$

$\square$

**Remark 2.2.8.** For non-negative random variables the above inequality can be improved as follows:
$$\mathbb{P}(X = 0) \leq \frac{\text{Var}(X)}{\mathbb{E}(X^2)}.$$

Since $\mathbb{E}(X^2) \geq (\mathbb{E}X)^2$ this is indeed an improvement. The proof of this inequality is a simple application of the Cauchy–Schwarz inequality: set $A = \{\omega \mid X(\omega) > 0\}$, then

$$\left(\int_A X dP\right)^2 \leq \left(\int_A 1 dP\right)\left(\int_A X^2 dP\right),$$

25

that is
$$(\mathbb{E}X)^2 \leq (1 - \mathbb{P}(X = 0))(\mathbb{E}(X^2)).$$

After some algebraic manipulation we get that
$$\mathbb{P}(X = 0) \leq \frac{\mathrm{Var}(X)}{\mathbb{E}(X^2)}.$$

Theorem 2.2.7 implies that if $\lim_{n\to\infty} \frac{\mathrm{Var}(X_n)}{(\mathbb{E}X_n)^2} = 0$ then
$$\lim_{n\to\infty} \mathbb{P}(X_n = 0) = 0.$$

We can also see that if $\mathrm{Var}(X_n) = o((\mathbb{E}X_n)^2)$ then $X_n$ is concentrated around $\mathbb{E}X_n$ which we can simply denote by $X_n \sim \mathbb{E}X_n$.

## 2.3 Existence results

In this section we give the most basic examples of the probabilistic method where one only needs to use the union bound, Proposition 2.2.3.

### 2.3.1 Diagonal Ramsey numbers

Recall that the Ramsey-number $R(r, b)$ denotes the smallest $n$ such that no matter how we color the edges of the complete graph $K_n$ with red and blue colors it will either contain an induced red $K_r$ or a blue $K_b$. Note that the definition implies that for $n = R(r, b) - 1$ there is a coloring of $K_n$ without red $K_r$ and blue $K_b$.

**Theorem 2.3.1** (Erdős). *Suppose that the positive integers $n, k$ satisfy the inequality $\binom{n}{k} 2^{1-\binom{k}{2}} < 1$. Then $R(k, k) > n$. In particular, $R(k, k) > \lfloor 2^{k/2} \rfloor$ if $k \geq 3$.*

*Proof.* We need to show that there exists a coloring of the edge set of $K_n$ that does not contain either monochromatic red or blue clique $K_k$. Let us color each edges with color red or blue with probability $1/2$ independently of each other. Now let us estimate the probability that the coloring is bad, i. e., it contains a monochromatic red or blue $K_k$. For each $S \subset V(G)$ with $|S| = k$ let $A_S$ be the event the induced subgraph on $S$ is monochromatic. Then
$$\mathbb{P}(\text{coloring is bad}) \leq \sum_{|S|=k} \mathbb{P}(A_S) = \binom{n}{k} \frac{2}{2^{\binom{k}{2}}}.$$

By the condition of the theorem $\binom{n}{k}2^{1-\binom{k}{2}} < 1$, so the probability that the coloring is good is positive.

Next we show that for $k \geq 3$ and $n = \lfloor 2^{k/2} \rfloor$ the condition of the theorem is satisfied. Indeed,

$$\binom{n}{k}2^{1-\binom{k}{2}} < \frac{n^k}{k!}2^{1-\binom{k}{2}} \leq \frac{2^{k^2/2}}{k!}2^{1-\binom{k}{2}} = \frac{2^{(k+2)/2}}{k!} < 1.$$

if $k \geq 3$. $\qquad\square$

### 2.3.2 Tournaments

**Definition 2.3.2.** A tournament is a complete directed graph. A tournament $D$ is called $k$-dominated if for every $k$ vertices $v_1, \ldots, v_k$ there exists a vertex $u$ such that $(u, v_i) \in E(D)$ for $i = 1, \ldots, k$.

**Theorem 2.3.3** (Erdős [6])**.** *If $n$ is large enough then there exists a $k$-dominated tournament on $n$ vertices.*

*Proof.* Let us direct each edge with probability $1/2 - 1/2$ independently of each other. Then the chance that for a given set of vertices $v_1, \ldots, v_k$ there is no $u$ such that all $(u, v_i) \in E(D)$ is $(1 - 1/2^k)^{n-k}$. Hence the probability that the orientation is bad is at most

$$\binom{n}{k}\left(1 - \frac{1}{2^k}\right)^{n-k}.$$

A little computation shows that if $\frac{n}{\ln n} > k2^k$ then this is less than 1. For large $k$ this is satisfied if $n > k^2 2^k$. Hence with positive probability there exists a $k$-dominated tournament. $\qquad\square$

**Remark 2.3.4.** The computation can be carried out using the bound $1 + x < e^x$:

$$\begin{aligned}
\binom{n}{k}\left(1 - \frac{1}{2^k}\right)^{n-k} &\leq \frac{n^k}{k!}\exp\left(-\frac{1}{2^k}(n-k)\right) \\
&\leq \frac{1}{k!}\exp\left(\frac{k}{2^k}\right)\cdot\exp\left(k\ln n - \frac{n}{2^k}\right) \\
&\leq \frac{e}{k!}\exp\left(k\ln n - \frac{n}{2^k}\right) \\
&\leq \frac{e}{k!} < 1
\end{aligned}$$

if $\frac{n}{\ln n} > k2^k$.

## 2.4 First moment method

In the previous section we have seen some very simple ideas how to find a certain structure $S$ by proving that it exists with positive probability just by using union bound. Here we study another very simple technique. This is the so-called first moment method. In many cases the structure $S$ that we need to find is defined through some parameter $f(S)$. For instance, we need to prove that there exists a structure $S$ for which some parameter $f(S)$ satisfies $f(S) \geq \rho$. If we find a probability space in which the expected value of $f(S)$ is bigger or equal to $\rho$ then we can immediately conclude that $f(S) \geq \rho$ with positive probability.

### 2.4.1 Warm-up: large bipartite subgraphs

**Theorem 2.4.1** ([2]). *Let $G$ be a graph with $n$ vertices and $e(G)$ edges. Then $G$ has a bipartite subgraph with at least $e(G)/2$ edges.*

*Proof.* One can rephrase the statement of the theorem as follows: there exists a cut $(A, V \setminus A)$ of $G$ such that the number of edges $(e(A, V \setminus A))$ contained in the cut is at least $e(G)/2$.

Let us consider the random set $A$ which contains every $v \in V(G)$ with probability $1/2$ independently of each other. (This way we have defined a probability space.) Let us consider the random variable $X = e(A, V \setminus A)$. We have to show that with positive probability $X \geq e(G)/2$. To this end it is enough to show that $\mathbb{E}X = e(G)/2$. This is indeed true. For every edge $f \in E(G)$ let us introduce the indicator random variable $X_f$ which takes value $1$ if $f$ is in the cut $(A, V \setminus A)$, and $0$ otherwise. Then

$$\mathbb{E}X = \mathbb{E}\left( \sum_{f \in E(G)} X_f \right) = \sum_{f \in E(G)} \mathbb{E}X_f.$$

(Note that the random variables $X_f$ are not necessarily independent, but the linearity of expectation holds true even with non-independent random variables.) For all $f \in E(G)$ we have $\mathbb{E}X_f = 1/2$ since the end points of $f$ are in the same set with probability $1/2$ and they are in different sets with probability $1/2$. Hence

$$\mathbb{E}X = \sum_{f \in E(G)} \mathbb{E}X_f = \sum_{f \in E(G)} \frac{1}{2} = \frac{1}{2}e(G).$$

We are done! $\qquad \square$

## 2.4.2 Independent sets

**Theorem 2.4.2** (Caro; Wei)**.** *Let $G$ be a graph with vertex degrees $d_1, \ldots, d_n$. Let $\alpha(G)$ be the size of the largest independent set of the graph $G$. Then*

$$\alpha(G) \geq \sum_{i=1}^{n} \frac{1}{d_i + 1}.$$

*Proof.* Consider a random permutation of the vertices. Let us encircle all the vertices that precede all their neighbors in the given order. Let $X(\pi)$ be the random variable that counts the number of encircled vertices. For a given vertex $v \in V(G)$ let $X_v$ be the indicator variable that the the vertex $v$ is encircled or not. Then $X = \sum_{v \in V(G)} X_v$, consequently

$$\mathbb{E}X = \sum_{v \in V(G)} \mathbb{E}X_v.$$

Note that for a vertex $v$ we have $\mathbb{E}X_v = \frac{1}{d_v+1}$ since the probability that $v$ precedes its neighbors is the same as saying that $v$ is the first among $d_v + 1$ vertices in a random permutation, and this probability is clearly $\frac{1}{d_v+1}$. Hence

$$\mathbb{E}X = \sum_{v \in V(G)} \mathbb{E}X_v = \sum_{i=1}^{n} \frac{1}{d_i + 1}.$$

With positive probability $X$ is at least as large as this expected value. On the other hand, in an arbitrary order the encircled vertices form an independent set since if two of them were adjacent then the second of the two vertices in the order would not be encircled. Hence

$$\alpha(G) \geq \mathbb{E}X = \sum_{i=1}^{n} \frac{1}{d_i + 1}$$

as required. $\qquad\square$

**Remark 2.4.3.** From the above proof one can easily deduce Turán's theorem.

## 2.4.3 Crossing number

**Theorem 2.4.4** (Ajtai-Chvátal-Newborn-Szemerédi [1]; Leighton)**.** *Let $G$ be a graph with $n$ vertices and $e$ edges. Let $X(G)$ be the crossing number of the graph $G$. If $e \geq 4n$ then*

$$X(G) \geq \frac{e^3}{64n^2}.$$

*Proof.* Recall that any planar graph with $n$ vertices has at most $3n - 6$ edges. Consequently, if $G$ is a graph with $n$ vertices and $e$ edges then the crossing number is at least $e - (3n - 6)$ (why?). So

$$X(G) \geq e(G) - 3v(G).$$

(The +6 won't be important for us.) This is of course a weaker statement than what we want to prove. The key idea of the better bound is to apply this weak inequality to a random subgraph of $G$. Set $0 \leq p \leq 1$ and consider the random subgraph of $G$ where we keep each vertex with probability $p$ and delete it with probability $1 - p$. Let $G_p$ be the obtained graph. Then

$$\mathbb{E}v(G_p) = pv(G) \quad \text{and} \quad \mathbb{E}e(G_p) = p^2 e(G),$$

since the probability that we keep an edge is $p^2$, the probability that we keep both end points of the edge. We need to be a bit more careful with $\mathbb{E}X(G_p)$. Starting from an optimal drawing of $G$, the probability that a crossing remains is $p^4$ since all four vertices determining the crossing should remain. This means that starting from an optimal drawing of $G$ the expected value of the crossing number of $G_p$ is $p^4 X(G)$. However, it may happen that $G_p$ has a better drawing with smaller number of crossings. So all we can say is that

$$\mathbb{E}X(G_p) \leq p^4 X(G).$$

Hence
$$p^4 X(G) - p^2 e(G) + 3pv(G) \geq \mathbb{E}X(G_p) - \mathbb{E}e(G_p) + 3\mathbb{E}v(G_p) =$$
$$= \mathbb{E}(X(G_p) - e(G_p) + 3v(G_p)) \geq 0.$$

Whence $p^4 X(G) - p^2 e(G) + 3pv(G) \geq 0$ for all $0 \leq p \leq 1$. Now let us choose $p$ to be $\frac{4v(G)}{e(G)}$. This is at most 1 according to the assumption of the theorem. Then

$$X(G) \geq p^{-2} e(G) - 3p^{-3} v(G) = \frac{e(G)^3}{64v(G)^2}.$$

This is exactly what we wanted to prove. $\qquad\square$

Since it is not clear how we can use such a statement let us consider a corollary of this theorem. Then later we even consider a corollary of this corollary.

Given some points and lines on the plane. Let $\mathcal{P}$ be the set of points, and $\mathcal{L}$ be the set of lines. The number of point-line incidences is exactly what we expect:

$$I(\mathcal{P}, \mathcal{L}) = |\{(P, L) \in \mathcal{P} \times \mathcal{L} \mid P \in L\}|.$$

Let $I(n, m)$ be the maximal number of incidences given $n$ points and $m$ lines:

$$I(n, m) = \max_{|\mathcal{P}|=n, |\mathcal{L}|=m} I(\mathcal{P}, \mathcal{L}).$$

The following theorem gives a good bound on $I(n, m)$.

**Theorem 2.4.5** (Szemerédi-Trotter [11])**.**

$$I(n, m) \leq 4(m^{2/3}n^{2/3} + m + n).$$

*Proof.* Let us consider the graph $G$ whose vertices are the elements of the set $\mathcal{P}$, i. e., the points, and two points are adjacent if there is a line $\ell \in \mathcal{L}$ that contains the two points next to each other.

First let us determine the number of edges of the graph $G$. If a line contains $k$ points then it determines $k - 1$ edges. Hence the number of edges is $I(\mathcal{P}, \mathcal{L}) - m$. Next let us give an upper bound on $X(G)$. Two edges intersect each other if two lines intersect each other. Hence $X(G)$ is at most $\binom{m}{2}$. If $e(G) < 4n$ then

$$I(\mathcal{P}, \mathcal{L}) < 4n + m < 4(m^{2/3}n^{2/3} + m + n).$$

If $e(G) \geq 4n$ then we can use the previous theorem:

$$\binom{m}{2} \geq X(G) \geq \frac{e(G)^3}{64n^2} = \frac{(I(\mathcal{P}, \mathcal{L}) - m)^3}{64n^2}.$$

Thus

$$I(\mathcal{P}, \mathcal{L}) \leq (32m^2n^2)^{1/3} + m < 4(m^{2/3}n^{2/3} + m + n).$$

Hence

$$I(n, m) \leq 4(m^{2/3}n^{2/3} + m + n).$$

$\square$

**Remark 2.4.6.** We used very little information about the lines. We simply used that two lines have at most one intersection. We could have considered circles or arbitrary curves of degree at most $d$, these curves have also bounded number of intersections. Naturally, the constants in the theorem would have been worse, but still we would have received a bound of type $O_d(n^{2/3}m^{2/3} + n + m)$ for the number of incidences.

In what follows we consider a nice application of the Szemerédi-Trotter theorem. Let $A \subset \mathbb{R}$ be a finite set, and let

$$A + A = \{a + a' \mid a, a' \in A\}$$

and

$$A \cdot A = \{a \cdot a' \mid a, a' \in A\}.$$

If $A = \{1, 2, \ldots, n\}$ then $A+A = \{2, \ldots, 2n\}$, and so $|A+A| = 2n-1$. However, in this case we have $|A \cdot A| = \Omega\left(\frac{n^2}{(\log n)^\alpha}\right)$. If $A = \{1, 2, 2^2, \ldots, 2^{n-1}\}$ then $|A \cdot A| = 2n-1$, but then we have $|A + A| = \binom{n}{2}$. After checking several examples one will have the feeling that one of the sets should be large. This is a well-known conjecture:

**Conjecture 2.4.7** (Erdős-Szemerédi). For all $\varepsilon > 0$ there exists a constant $c(\varepsilon)$ such that for all finite set $A \subset \mathbb{R}$ we have

$$|A + A| + |A \cdot A| \geq c(\varepsilon)|A|^{2-\varepsilon}.$$

We are very far from proving this conjecture. The following result of György Elekes was a real breakthrough in 1997, and it opened the way of geometric arguments in additive combinatorics.

**Theorem 2.4.8** (Elekes [4]). *Let $A \subset \mathbb{R}$ be a finite set. Then*

$$|A + A| \cdot |A \cdot A| \geq c|A|^{5/2}.$$

*In particular,*

$$|A + A| + |A \cdot A| \geq c'|A|^{5/4}.$$

*Proof.* Let $P = \{(a, b) \mid a \in A + A, b \in A \cdot A\}$. This is a point set on the plane of size $|A + A||A \cdot A|$.

Let us consider the lines of following type:

$$\ell_{a,b} = \{(x, y) \mid y = a(x - b)\},$$

where $a, b \in A$. Let $L$ be the set of these lines. Then $|L| = |A|^2$. Every such line contains $|A|$ points form $P$: $(b + c, ac) \in \ell_{a,b}$ if $c \in A$. Whence $I(P, L) \geq |A|^3$. According to Szemerédi-Trotter theorem we have

$$|A|^3 \leq 4((|A + A| \cdot |A \cdot A|)^{2/3}(|A|^2)^{2/3} + |A + A| \cdot |A \cdot A| + |A|^2).$$

From this the statement of the theorem follows after a little computation. $\qquad\square$

**Remark 2.4.9.** Currently, the best-known result is due to József Solymosi [9]:

$$|A + A| + |A \cdot A| \geq c(\varepsilon)|A|^{4/3-\varepsilon}.$$

More precisely, Solymosi showed that

$$|A \cdot A| \cdot |A + A|^2 \geq \frac{|A|^4}{4\lceil \ln |A| \rceil},$$

consequently,

$$\max(|A \cdot A|, |A + A|) \geq \frac{|A|^{4/3}}{2\lceil \ln |A| \rceil^{1/3}}.$$

## 2.5 Alteration

In this section we study a method called the altered first moment method. It is a slightly bit more tricky than the first moment method. Here the randomly chosen structure $S$ will not be immediately good, but will be bad just a little bit so that we can fix the bad part of the structure. In practice, there will be a parameter $f(.)$ that measures the badness of the structure (or if there is a given parameter $f(.)$ already, then we prepare a new parameter $f'(.)$ measuring $f(.)$ and the badness at the same time). If the expected value of this badness parameter is small, then with positive probability we can find a random structure that we can fix later. After the examples it will be clear how this method works.

### 2.5.1 Independent sets in graphs and hypergraphs

**Theorem 2.5.1.** *Let $H$ be an $r$-uniform hypergraph with $n$ vertices and $e(H)$ edges. Suppose that $n \leq 2e$. Then there exists a set $S \subseteq V(H)$ inducing no edge such that*

$$|S| \geq \frac{1}{2}\left(\frac{n}{2e(H)}\right)^{1/(r-1)} n.$$

*Proof.* Let $T$ be a random subset of the vertex set chosen as follows: we choose each element of $V$ to be in $T$ with probability $p$. We will choose $p$ later. Then $\mathbb{E}|T| = pn$ and for the number of edges induced by $T$ we have $\mathbb{E}(e(T)) = p^r e(H)$. Then

$$\mathbb{E}(|T| - e(T)) = pn - p^r e(H).$$

Let

$$p = p_0 = \left(\frac{n}{2e(H)}\right)^{1/(r-1)}.$$

Then $p_0 n - p_0^r e(H) = p_0 n/2$. Therefore,

$$\mathbb{E}(|T| - e(T)) = \frac{1}{2}p_0 n.$$

Hence there must be a set $T$ for which $|T| - e(T) \geq \frac{1}{2}p_0 n$. Let $S \subseteq T$ be a set obtained from $T$ by deleting one vertex of each edge of $T$. Then $S$ induces no edge and

$$|S| \geq |T| - e(T) \geq \frac{1}{2}p_0 n = \frac{1}{2}\left(\frac{n}{2e(H)}\right)^{1/(r-1)} n.$$

$\square$

**Remark 2.5.2.** In the case of graphs, that is $r = 2$, this theorem says that

$$\alpha(G) \geq \frac{n^2}{4e(G)}.$$

This is always weaker than the bound

$$\alpha(G) \geq \sum_{i=1}^{n} \frac{1}{d_i + 1}$$

obtained earlier.

We could have chosen $p$ in a bit better way by simply choosing it such a way that it maximizes $pn - p^r e(H)$. This would have yielded the bound

$$\alpha(H) \geq \frac{r-1}{r}\left(\frac{n}{re(H)}\right)^{1/(r-1)} n.$$

### 2.5.2 Ramsey-numbers revisited

**Theorem 2.5.3** ([2]). *For all $n$ and $k$ we have $R(k,k) > n - \binom{n}{k}2^{1-\binom{k}{2}}$.*

*Proof.* Let us color the edges of a complete graph $K_n$ with red and blue. Let $X$ be the number of monochromatic $K_k$. Then

$$\mathbb{E}X = \binom{n}{k}2^{1-\binom{k}{2}}.$$

So there must be a coloring with at most as many monochromatic $K_k$. Now let us delete one vertex from each monochromatic $K_k$. Then the number of vertices is at least $n - \binom{n}{k}2^{1-\binom{k}{2}}$ and the resulting graph has no monochromatic $K_k$. Hence $R(k,k) > n - \binom{n}{k}2^{1-\binom{k}{2}}$. $\square$

**Remark 2.5.4.** A careful analysis shows that this bound implies that

$$R(k,k) \geq \frac{1}{e}(1 + o(1))k2^{k/2}$$

while our previous argument only gave

$$R(k,k) \geq \frac{1}{\sqrt{2}e}(1 + o(1))k2^{k/2}.$$

Further improvement can be obtained by the so-called Lovász local lemma:

$$R(k,k) \geq \frac{\sqrt{2}}{e}(1 + o(1))k2^{k/2}.$$

### 2.5.3  Dominating sets in graphs

**Theorem 2.5.5** ([2]). *Let $G = (V, E)$ be a graph with $n$ vertices and minimum degree $\delta > 1$. Then it has a dominating set of size at most*

$$n\frac{1 + \ln(\delta + 1)}{\delta + 1}.$$

*(A set $U$ is called a dominating set of $G$ if all $v \in V \setminus U$ has some neighbor $u$ in $U$.)*

*Proof.* The strategy is the following: we choose a random subset $S$ and let $T = T(S)$ be the set of vertices $v$ such that neither $v$, nor any of the neighbors of $v$ are in the set $S$. Then $S \cup T$ is a dominating set. Let us choose $S$ as follows: we choose each vertex $v$ into $S$ with probability $p$. We will choose $p$ later. Then for any vertex $v \in V$ we have

$$\mathbb{P}(v \in T) = (1 - p)^{1+d(v)} \leq (1 - p)^{1+\delta} \leq e^{-p(\delta+1)}$$

since neither $v$, nor any of the neighbors of $v$ are in the set $S$. Hence

$$\mathbb{E}(|S| + |T|) = \mathbb{E}|S| + \mathbb{E}|T| \leq n(p + e^{-p(\delta+1)}).$$

Let

$$p = \frac{\ln(\delta + 1)}{\delta + 1}.$$

Then

$$\mathbb{E}(|S| + |T|) \leq n(p + e^{-p(\delta+1)}) = \frac{n(1 + \ln(\delta + 1))}{\delta + 1}.$$

Hence with positive probability there must be a dominating set of at most this size. $\qquad\square$

### 2.5.4 Graphs with large chromatic number and girth

**Theorem 2.5.6** (Erdős [5]). *For arbitrary $(k, \ell)$ there exists a graph $G$ whose chromatic number is at least $k$ and the length of its shortest cycle is at least $\ell$.*

*Proof.* Let $G(n, p)$ be the random graph with $n$ vertices such that we draw all edges with probability $p = p(n)$ independently of each other. In this proof we will set $p = n^{-\alpha}$, where $\alpha \geq 0$ is a parameter chosen later. First we estimate the number of cycles shorter than $\ell$ . Given vertices $v_1 v_2 \ldots v_r$ form a cycle if $v_i v_{i+1}$ $(r + 1 = 1)$ are all edges, the probability of this event is $p^r$. Naturally, we can choose the sequence $v_1 v_2 \ldots v_r$ in $n(n-1) \ldots (n-r+1)$ ways, we only have to take take into account that we counted the same cycle $2r$ ways (rotated and reflected copies). Let $X$ be the random variable counting the number of cycles of length at most $\ell - 1$. Furthermore, let $X(v_1 \ldots v_r)$ $(r \leq \ell - 1)$ be the indicator random variable that the vertices $v_1 \ldots v_r$ form a cycle in this order. Then

$$X = \sum_{r, v_1 \ldots v_r} X(v_1 \ldots v_r).$$

Hence

$$\mathbb{E}X = \sum_{r, v_1 \ldots v_r} \mathbb{E}X(v_1 \ldots v_r) = \sum_{r=3}^{\ell-1} \frac{n(n-1) \ldots (n-r+1)}{2r} p^r \leq \sum_{r=3}^{\ell-1} \frac{(np)^r}{2r}.$$

Set $M = \sum_{r=3}^{\ell-1} \frac{(np)^r}{2r}$. Suppose that with some choice of $p$ we can ensure that $M$ is small then with positive probability the number of cycles of length at most $\ell - 1$ will be at most $M$ and by throwing out one point from each cycle we get a graph on at least $n - M$ vertices that does not contain a cycle of length at most $\ell - 1$. In fact, we need to be a little bit more careful as we need that the number of short cycles is small with large probability. Fortunately, we get it immediately: with probability at least $1/2$ the number of cycles of length at most $\ell - 1$ is at most $2M$. Otherwise the expected value would be bigger than $M$.

Before we try to chose $p$ appropriately let us see how we can bound the chromatic number $\chi(G)$ of $G$. Here we use the simple fact that

$$\chi(G) \geq \frac{n}{\alpha(G)}.$$

This is true since all coloring class induces an independent set so its size is at most $\alpha(G)$, so we need at least $\frac{n}{\alpha(G)}$ colors to color $G$. So to make $\chi(G)$ large, it is enough

to ensure that $\alpha(G)$ is small. Let us bound the probability that $\alpha(G) \geq s$. For a set $S$ of size $s$ let $A_S$ be the event that $S$ does not induce any edge. Then

$$\mathbb{P}(\alpha(G) \geq s) \leq \sum_{|S|=s} \mathbb{P}(A_S) = \binom{n}{s}(1-p)^{\binom{s}{2}} \leq n^s(1-p)^{\binom{s}{2}} \leq (ne^{-p(s-1)/2})^s.$$

(In the last step we used the fact that $1 + x \leq e^x$ is satisfied for all $x$. This is a rather standard bound that is quite good if $x$ is small.)

Now it is clear what we have to keep in mind: let $M$ be small, so we need a small $p$, but we also need that $s$ is not too large and so we need that $ne^{p(s-1)/2} < 1$. We can easily achieve it as follows: set $p = n^{\theta-1}$ where $\theta = \frac{1}{2(\ell-1)}$ and $s = \lceil \frac{3}{p} \log n \rceil$. Then

$$M = \sum_{r=3}^{\ell-1} \frac{(np)^r}{2r} \leq n^{\theta(\ell-1)} \sum_{r=3}^{\ell-1} \frac{1}{2r} \leq n^{1/2} \log n \leq \frac{n}{4}$$

if $n$ is large enough. On the other hand,

$$\mathbb{P}(\alpha(G) \geq s) \leq (ne^{-p(s-1)/2})^s \leq 1/4$$

if $n$ is large enough. Since $\mathbb{P}(X \geq 2M) \leq 1/2$ and $\mathbb{P}(\alpha(G) \geq s) \leq 1/4$, with positive probability there exists a graph where the number of short cycles is at most $n/2$ and $\alpha(G) \leq s$. Now from all cycles of length at most $\ell - 1$ let us throw out 1 vertex and let $G^*$ be the obtained graph. Then $G^*$ has at least $n/2$ vertices and it does not contain a cycle of length at most $\ell - 1$. Furthermore, $\alpha(G^*) \leq \alpha(G)$ since $G^*$ is an induced subgraph of $G$. Then

$$\chi(G^*) \geq \frac{|V(G^*)|}{\alpha(G^*)} \geq \frac{n/2}{3n^{1-\theta} \log n} = \frac{n^\theta}{6 \log n}.$$

If $n$ is large enough this is bigger than $k$. We are done! $\qquad\square$

## 2.6   Second Moment method

In the previous sections we used the union bounds and the first moment method. These techniques are very powerful due to the fact that they do not require any information on the dependence of the random variables.

In this section we see some applications of the second moment method which roughly means that we use Chebyshev's inequality as a new ingredient in our proofs. We will see that at least we need some partial information about the dependence

of the random variables, but not too much. Generally quite crude bounds will be enough to achieve our goals.

This section is based on the corresponding chapter of the book The Probabilistic Method by Noga Alon and Joel Spencer.

In this section we study the threshold function of random graphs. This topic was initiated in the seminal paper [7] of Erdős and Rényi: *On the evolution of random graphs*, Magyar Tud. Akad. Mat. Kutató Int. Közl. 5 (1960), 17-61. In fact, all results of this section can be found in this paper. This paper is on the internet in a scanned form.

### 2.6.1  General approach

From Section 2.2 we know that for a non-negative random variable $X$ taking only integers, then we have

$$1 - \mathbb{E}X \leq \mathbb{P}(X = 0) \leq \frac{\mathrm{Var}(X)}{(\mathbb{E}X)^2}.$$

These two inequalities will play a major role in this section. We will often encounter the situation that having some property is equivalent with some random variable taking value 0. Hence if $\mathbb{E}X$ is small then $\mathbb{P}(X = 0)$ is large, and so the random structure has the desired property with large probability. On the other hand, if $\frac{\mathrm{Var}(X)}{(\mathbb{E}X)^2}$ is small then the random structure doesn't have the desired property with large probability.

Often we will encounter with a sequence of structures, notably a sequence of random graphs $G(n, p)$. In this case $X$ will be some $X_n$ in a sequence. As we will see it is also worth considering separately the case when $X_n = X_1^{(n)} + X_2^{(n)} + \cdots + X_m^{(n)}$ where $X_i^{(n)}$ are indicator random variables. Let $X_i^{(n)}$ be the indicator random variable of the event $A_i^{(n)}$. Let us introduce the notation $i \sim j$ if $A_i^{(n)}$ and $A_j^{(n)}$ are not independent. Then it is also worth introducing the following sum:

$$\Delta_n = \sum_{i \sim j} \mathbb{P}(A_i^{(n)} \cap A_j^{(n)}).$$

(In this sum both $(i, j)$ and $(j, i)$ appear.) If $\mathbb{P}(A_i^{(n)}) = p_i^{(n)}$ then

$$\mathrm{Var}(X_i^{(n)}) = \mathbb{E}(X_i^{(n)})^2 - (\mathbb{E}X_i^{(n)})^2 = p_i^{(n)} - (p_i^{(n)})^2 \leq p_i^{(n)} = \mathbb{E}X_i^{(n)}.$$

Furthermore,

$$\mathrm{Cov}(X_i^{(n)}, X_j^{(n)}) = \mathbb{E}(X_i^{(n)} X_j^{(n)}) - \mathbb{E}X_i^{(n)} \cdot \mathbb{E}X_j^{(n)} \leq \mathbb{E}(X_i^{(n)} X_j^{(n)}) = \mathbb{P}(A_i^{(n)} \cap A_j^{(n)}).$$

Using these inequalities we get that

$$\text{Var}(X_n) = \sum_{i=1}^{n} \text{Var}(X_i^{(n)}) + 2 \sum_{i<j} \text{Cov}(X_i^{(n)}, X_j^{(n)})$$

$$= \sum_{i=1}^{n} \text{Var}(X_i^{(n)}) + \sum_{i \sim j} \text{Cov}(X_i^{(n)}, X_j^{(n)})$$

$$\leq \sum_{i=1}^{n} \mathbb{E}X_i^{(n)} + \sum_{i \sim j} \mathbb{P}(A_i^{(n)} \cap A_j^{(n)})$$

$$= \mathbb{E}X_n + \Delta_n.$$

Here we used the fact that if $i \nsim j$, equivalently $A_i^{(n)}$ and $A_j^{(n)}$ are independent, then $\text{Cov}(X_i^{(n)}, X_j^{(n)}) = 0$. Hence

$$\text{Var}(X_n) \leq \mathbb{E}X_n + \Delta_n.$$

Hence Theorem 2.2.7 implies the following statement.

**Theorem 2.6.1.** *Suppose that $\mathbb{E}X_n \to \infty$ and $\Delta_n = o((\mathbb{E}X_n)^2)$.*
*Then $\mathbb{P}(X_n > 0) \to 1$.*

It is worth doing some extra work with $\Delta_n$. Many times the indicator random variables $X_1^{(n)}, \ldots, X_m^{(n)}$ have a symmetric role, in other words, for all $i$ and $j$ there is an automorphism of the underlying space that takes $A_i^{(n)}$ to $A_j^{(n)}$. Then

$$\Delta_n = \sum_{i \sim j} \mathbb{P}(A_i^{(n)} \cap A_j^{(n)}) = \sum_{i} \mathbb{P}(A_i^{(n)}) \sum_{j \sim i} \mathbb{P}(A_j^{(n)} \mid A_i^{(n)}).$$

The inner sum is independent of $i$, because of the symmetry:

$$\Delta_n^* = \sum_{j \sim i} \mathbb{P}(A_j^{(n)} \mid A_i^{(n)}).$$

Hence

$$\Delta_n = \sum_{i} \mathbb{P}(A_i^{(n)}) \Delta_n^* = \Delta_n^* \sum_{i} \mathbb{P}(A_i^{(n)}) = \Delta_n^* \mathbb{E}X_n.$$

So in this case we get the following theorem

**Theorem 2.6.2.** *Suppose that $\mathbb{E}X_n \to \infty$ and $\Delta_n^* = o(\mathbb{E}X_n)$. Then $\mathbb{P}(X_n > 0) \to 1$.*

**Remark 2.6.3. (Important!)** It is rather inconvenient to write out the $.^{(n)}$ every time: $X_i^{(n)}, A_i^{(n)}, p_i^{(n)}...$ So in what follows we hide the notation $n$ and for instance the last claim will read as follows: "Suppose that $\mathbb{E}X \to \infty$ and $\Delta^* = o(\mathbb{E}X)$. Then $\mathbb{P}(X > 0) \to 1$." This is of course completely stupid if we forget that there is a hidden parameter $n$. Nevertheless, the parameter $n$ will always be clear from the context. For instance, if we study the random graph $G(n, p(n))$ and $X$ is the number of $K_4$ in the graph then it is clear that actually $X = X_n$ belongs to $G(n, p(n))$.

## 2.6.2 Threshold functions of graph appearance

Let $G(n, p)$ be the random graph on $n$ vertices whose edges appear with probability $p$ independently of each other. The probability $p$ may depend on $n$, for instance, it can be $p = p(n) = n^{-1/2}$.

Surprisingly, one can see "all" graphs $G(n, p)$ at the same time as $p$ runs from 0 to 1. For all edges let us pick a random number from the interval $[0, 1]$, then just as we rotate the frequency finder of a radio we start to increase $p$. At some point $t$ the edges with a number less than $t$ will lit up. As we increase $t$ more and more edges will lit up. At point $t = 0$ the whole graph is dark (with probability 1, while at $t = 1$ the whole graph is lit up. At point $p$ we can see $G(n, p)$. This process is called the *evolution of random graphs*.

What kind of questions can we study? We can, for instance, ask for the probability that $G(n, p)$ contains a Hamiltonian-cycle or we can seek for the probability that the graph is a planar graph or the probability that its chromatic number is at most 100. For a fixed $n$ these questions might be very difficult to answer and answers might be very ugly. In general, we only wish to know the answer as the number of vertices tends to infinity. In other words, we are seeking $\lim \mathbb{P}(G(n, p) \in P)$ for some property $P$ like containing Hamiltonian-cycle or not. Actually, we will be even less ambitious as we only try to determine the so-called threshold function of the property $P$.

For a property $P$ a function $p_t(n)$ is the threshold function if

$$\lim_{n\to\infty} \mathbb{P}(G(n, p(n)) \in P) = \begin{cases} 1 & \text{if } \frac{p(n)}{p_t(n)} \to \infty, \\ 0 & \text{if } \frac{p(n)}{p_t(n)} \to 0. \end{cases}$$

If the probability of a (sequence of) events converge to 1 then we simply say that the considered event asymptotically almost surely happens. From the definition of

the threshold function it is clear that being a threshold function is not a uniquely determined function. For instance, if $p_t(n)$ is a threshold function then for any positive constant $c$ the function $cp_t(n)$ is also a threshold function. Another observation is that the definition suggests that we only consider a threshold function if increasing $p(n)$ also increases $\mathbb{P}(G(n, p(n)) \in P)$. This happens if the property $P$ is monotone increasing, this means that if $G$ has property $P$, then adding edges to $G$ won't lead out from $P$. For instance, if $G$ has a Hamiltonian-cycle and we add some edges, then the obtained graph will have a Hamiltonian-cycle too. If the chromatic number is at least 100, then no matter how many edges we add the chromatic number will be at least 100. But for instance, if we consider the planarity of $G$ then we should study the property that at which $p$ will $G(n, p)$ likely to loose the planarity. This is a monotone decreasing property.

$$\star \quad \star \quad \star$$

Now let us consider a concrete example: at which $p$ the graph $K_4$ will appear in $G(n, p)$?

**Theorem 2.6.4.** *The threshold function of the appearance of $K_4$ is $n^{-2/3}$.*

**Remark 2.6.5.** We can rephrase the claim as follows: the threshold function of the property $\omega(G) \geq 4$ is $n^{-2/3}$.

*Proof.* Let $S$ be a subset of size 4 of $V(G)$, where $G = G(n, p)$ is a random graph. Let $A_S$ be the event that $S$ induces a $K_4$ in $G$, and let $X_S$ be the indicator random variable of $A_S$. Let $X$ be the number of $K_4$ in $G$. Then

$$X = \sum_{\substack{S \subseteq V(G) \\ |S|=4}} X_S.$$

Whence

$$\mathbb{E}X = \sum_{\substack{S \subseteq V(G) \\ |S|=4}} \mathbb{E}X_S = \binom{n}{4} p^6 \leq \frac{(pn^{2/3})^6}{24}.$$

If $p(n)n^{2/3} \to_{n \to \infty} 0$ then $\mathbb{E}X \to 0$ as $n \to \infty$. Hence using $\mathbb{P}(X = 0) \geq 1 - \mathbb{E}X$ we get that

$$\lim_{n \to \infty} \mathbb{P}(\omega(G) \geq 4) = 0.$$

Now suppose that $p(n)n^{2/3} \to_{n \to \infty} \infty$. Then $\mathbb{E}X \to \infty$ as $n \to \infty$. We will use Theorem 2.6.1; since all sets of size 4 looks the same way, the random variables $X_S$ are symmetric. Note that $S \sim T$ if $|S \cap T| \geq 2$, otherwise the events $A_S$ and $A_T$ are independent since they don't have a common edge. Let us fix a set $S$. Then then there are $6\binom{n-2}{2} = O(n^2)$ sets $T$ that intersects $S$ in 2 elements and there are $4\binom{n-3}{1} = O(n)$ sets $T$ intersecting $S$ in 3 vertices. In the former case we have $\mathbb{P}(A_T|A_S) = p^5$, while in the latter case $\mathbb{P}(A_T|A_S) = p^3$. Then

$$\Delta^* = O(n^2 p^5) + O(np^3) = o(n^4 p^6) = o(\mathbb{E}X)$$

since $p(n)n^{-2/3} \to \infty$. Hence by Theorem 2.6.1 the graph $K_4$ appears asymptotically almost surely. $\square$

Now let us consider the bit more general problem of determining the threshold function of the appearance of a given graph $H$. After a quick check of the proof concerning $K_4$ we see that the value $2/3$ comes from the ratio of the vertices and edges of $K_4$. This may prompt us to believe that this is also the answer for the general question, i. e., for any graph $H$ the threshold function is $n^{-v(H)/e(H)}$. There is a minor problem with this idea: in order to make sure that $H$ appears one needs that all subgraphs $H'$ also appears, and it might very easily occur that for some $H'$ the value $n^{-v(H')/e(H')}$ is bigger than $n^{-v(H)/e(H)}$. This motivates the following definition.

**Definition 2.6.6.** Let $H$ be a graph with $v$ vertices and $e$ edges. We call the quantity $\rho(H) = \frac{e}{v}$ the density of $H$. We say that a graph $H$ is balanced if for all subgraph $H'$ we have $\rho(H') \leq \rho(H)$. The graph $H$ is said to be strictly balanced if for all proper subgraph $H'$ we have $\rho(H') < \rho(H)$.

The proof of the next theorem practically does not require any new idea.

**Theorem 2.6.7.** *Let $H$ be a balanced graph with $n$ vertices and $e$ edges. Let $P_H$ be the property that $H$ is a (not necessarily induced) subgraph of a graph $G$. Then the threshold function of $P_H$ is $p = n^{-v/e}$.*

*Proof.* For all subsets $S$ of size $v$ let $A_S$ be the event that $H$ is a subgraph of $G[S]$. Then

$$p^e \leq \mathbb{P}(A_S) \leq v!p^e.$$

Let $X_S$ be the indicator random variable of $A_S$. Furthermore, set $X = \sum_{|S|=v} X_S$. Hence the event that $G$ contains $H$ occurs if and only if $X > 0$. By the linearity of expectation we get that

$$\mathbb{E}X = \sum_{|S|=v} \mathbb{E}X_S = \binom{n}{v} \mathbb{P}(A_S) = \Theta(n^v p^e).$$

Hence if $p(n)n^{e/v} \to 0$, then $\mathbb{E}X = o(1)$, thus $X = 0$ asymptotically almost surely.

Now suppose that $p(n)n^{e/v} \to \infty$. Then $\mathbb{E}X \to \infty$. Let us consider $\Delta^*$. (We can do it as the events $A_S$ are symmetric.). If $S \sim T$ then $2 \leq |S \cap T| \leq v - 1$. Then

$$\Delta^* = \sum_{T \sim S} \mathbb{P}(A_T|A_S) = \sum_{i=2}^{v} \sum_{|T \cap S|=i} \mathbb{P}(A_T|A_S).$$

Let $i$ be fixed. Then there are $\binom{v}{i}\binom{n-v}{v-i} = O(n^{v-i})$ ways to choose a set $T$ intersecting $S$ in exactly $i$ vertices. The subgraph induced by $S \cap T$ has $i$ vertices and since $H$ was balanced, the intersection contains at most $i\frac{e}{v}$ edges. So there are at least $e - i\frac{e}{v}$ edges of $T$ not in the intersection with $S$. Whence

$$\mathbb{P}(A_T|A_S) = O(p^{e-i\frac{e}{v}}).$$

Hence

$$\Delta^* = \sum_{i=2}^{v-1} O(n^{v-i}p^{e-i\frac{e}{v}}) = \sum_{i=2}^{v-1} O((n^v p^e)^{1-i/v}) = \sum_{i=2}^{v-1} o(n^v p^e) = o(\mathbb{E}X)$$

since $n^v p^e \to \infty$. By Theorem 2.6.1 we get that $H$ appears in $G$ asymptotically almost surely. $\qquad\square$

Next we study the isolated vertices and connectedness of $G(n,p)$.

**Theorem 2.6.8.** *Let $\omega(n) \to \infty$. Furthermore, let $p_\ell(n) = (\log n - \omega(n))/n$ and $p_u(n) = (\log n + \omega(n))/n$. Then $G(n, p_\ell(n))$ contains an isolated vertex asymptotically almost surely while $G(n, p_u(n))$ does not contain isolated vertex vertex asymptotically almost surely.*

*Proof.* First we prove that $G(n, p_u(n))$ does not contain an isolated vertex asymptotically almost surely. From now on let $p = p_u(n)$. Let $X$ be the number of isolated vertices, and $X_v$ be the indicator random variable of $v$ being an isolated vertex. Then

$$X = \sum_{v \in V} X_v.$$

Observe that $\mathbb{P}(X_v = 1) = (1-p)^{n-1}$. We can assume that $p \leq 1/2$ (Why? Evolution!). Then

$$\mathbb{E}X = \sum_{v \in V} \mathbb{E}X_v = n(1-p)^{n-1} = \frac{1}{1-p}n(1-p)^n \leq 2ne^{-pn} = 2ne^{-\log n + \omega(n)} = 2e^{-\omega(n)} \to 0.$$

as $n \to \infty$. Then

$$\mathbb{P}(X = 0) \geq 1 - \mathbb{E}X \to 1.$$

Next we show that $G(n, p_\ell(n))$ contains an isolated vertex asymptotically almost surely. From now on let $p = p_\ell(n)$. As before, let $X$ be the number of isolated vertices, and $X_v$ be the indicator random variable of $v$ being an isolated vertex. Then $X = \sum_{v \in V} X_v$, and $\mathbb{P}(X_v = 1) = (1-p)^{n-1}$. Hence

$$\mathbb{E}X = \sum_{v \in V} \mathbb{E}X_v = n(1-p)^{n-1} \sim ne^{-\log n + \omega(n)} = e^{\omega(n)} \to \infty.$$

Let us determine $\mathbb{E}X^2$.

$$\mathbb{E}X^2 = \sum_{v \in V} \mathbb{E}X_v^2 + 2\sum_{u,v \in V} \mathbb{E}X_u X_v = n(1-p)^{n-1} + n(n-1)(1-p)^{2n-3}.$$

Whence

$$\mathrm{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2 = n(1-p)^{n-1} + n(n-1)(1-p)^{2n-3} - n^2(1-p)^{2(n-1)} \leq$$

$$\leq n(1-p)^{n-1} + n^2(1-p)^{2n-3} - n^2(1-p)^{2(n-1)} = n(1-p)^{n-1} + pn^2(1-p)^{2n-3} = \mathbb{E}X + \frac{p}{1-p}(\mathbb{E}X)^2$$

Thus

$$\mathbb{P}(X = 0) \leq \frac{\mathrm{Var}(X)}{(\mathbb{E}X)^2} \leq \frac{1}{\mathbb{E}X} + \frac{p}{1-p} \to 0$$

since

$$\frac{p}{1-p} \leq 2p \leq \frac{2\log n}{n},$$

if $n$ is large enough. Hence $G(n, p_a(n))$ contains an isolated vertex asymptotically almost surely.

$\square$

**Theorem 2.6.9.** *Let $\omega(n) \to \infty$. Furthermore, let $p_\ell(n) = (\log n - \omega(n))/n$ and $p_u(n) = (\log n + \omega(n))/n$. Then $G(n, p_\ell(n))$ is disconnected asymptotically almost surely while $G(n, p_u(n))$ is connected asymptotically almost surely.*

*Proof.* It is clear from the previous theorem that $G(n, p_\ell(n))$ is disconnected asymptotically almost surely since it contains an isolated vertex with high probability. So we only need to prove that $G(n, p_u(n))$ is connected asymptotically almost surely. This is stronger than what we proved earlier, namely that it does not contain an isolated vertex. From now on let $p = p_u(n)$, and let $X_k$ denote the number of connected components of size $k$. Furthermore, let

$$X = \sum_{k=1}^{\lfloor n/2 \rfloor} X_k.$$

This is the number of connected components of size at most $\lfloor n/2 \rfloor$. Note that if $G$ is connected, then $X = 0$, and if $G$ is disconnected then $X \geq 1$ non-negative integer. Since $\mathbb{P}(X = 0) \geq 1 - \mathbb{E}X$ we only need to prove that $\mathbb{E}X \to 0$ as $n \to \infty$.

Let $f(k, p)$ be the probability that a random graph $G(k, p)$ is connected. For a set $S$ let $X_S$ be indicator random variable that the graph induced by the set $S$ is a connected component of $G(n, p)$. Then

$$\mathbb{E}X_S = \mathbb{P}(X_S = 1) = f(|S|, p)(1 - p)^{|S|(n - |S|)}$$

since there must be no edge between $S$ and $V \setminus S$ and the induced subgraph must be connected. Then

$$\mathbb{E}X = \mathbb{E}\left( \sum_{1 \leq |S| \leq \lfloor n/2 \rfloor} X_S \right) = \sum_{k=1}^{\lfloor n/2 \rfloor} \binom{n}{k} f(k, p)(1 - p)^{k(n-k)}.$$

Since $f(k, p) \leq 1$ we have

$$\mathbb{E}X \leq \sum_{k=1}^{\lfloor n/2 \rfloor} \binom{n}{k} (1 - p)^{k(n-k)}.$$

We have

$$\sum_{k=1}^{\lfloor n/2 \rfloor} \binom{n}{k} (1 - p)^{k(n-k)} \leq \sum_{k=1}^{\lfloor n/2 \rfloor} \left( \frac{en}{k} \right)^k e^{-pk(n-k)}.$$

Here one term can be bounded as follows:

$$\left(\frac{en}{k}\right)^k e^{-pk(n-k)} = \exp\left(k(1 + \log n - \log k) - k(n-k)\frac{\log n + \omega(n)}{n}\right)$$

$$= \exp\left(-\omega(n)\frac{k(n-k)}{n}\right) \cdot \exp\left(k\left(1 + \frac{k}{n}\log n - \log k\right)\right)$$

$$\leq \exp\left(-\omega(n)\frac{n-1}{n}\right) \cdot \exp\left(k\left(1 + \frac{k}{n}\log n - \log k\right)\right)$$

$$\leq \exp\left(-\omega(n)\frac{n-1}{n}\right) e^{-k}.$$

if $300 \leq k \leq n/2$, and less than some constant $C\exp\left(-\omega(n)\frac{n-1}{n}\right)$ for $1 \leq k \leq 299$. Indeed, if $x = \frac{k}{n}$ then

$$2 + \frac{k}{n}\log n = 2 + x\log\frac{k}{x} = 2 + x\log\frac{1}{x} + x\log k \leq 2 + \frac{1}{2}\log 2 + \frac{1}{2}\log k \leq \log k$$

for $k \geq 300$. Then

$$\sum_{k=1}^{\lfloor n/2 \rfloor} \binom{n}{k}(1-p)^{k(n-k)} \leq \exp\left(-\omega(n)\frac{n-1}{n}\right) \cdot \left(\sum_{k=1}^{299} C + \sum_{k=300}^{\infty} e^{-k}\right) = C'\exp\left(-\omega(n)\frac{n-1}{n}\right).$$

This last expression goes to 0 as $n \to \infty$. Hence

$$\mathbb{P}(G(n,p) \text{ is not connected}) \to 0.$$

We are done. $\qquad\square$

**Remark 2.6.10.** Another way to estimate the sum

$$\sum_{k=1}^{\lfloor n/2 \rfloor} \binom{n}{k}(1-p)^{k(n-k)}$$

is the following.

$$\sum_{k=1}^{\lfloor n/2 \rfloor} \binom{n}{k}(1-p)^{k(n-k)} \leq \sum_{k=1}^{\lfloor n/2 \rfloor} \left(\frac{en}{k}\right)^k e^{-pk(n-k)} = \sum_{k=1}^{\lfloor n/2 \rfloor} \left(\frac{en}{k}e^{pk}e^{-pn}\right)^k = \sum_{k=1}^{\lfloor n/2 \rfloor} \left(e^{1-\omega(n)}\frac{e^{pk}}{k}\right)^k.$$

The function $e^{px}/x$ is convex on the interval $[0,\infty)$ for arbitrary $p$. In particular, it takes its maximum at one of the end points on the interval $[1, n/2]$. At 1 this function is at most $e$. At $n/2$ we can assume that $\omega(n) \leq \log n$ and we get that the value of the function is at most 2. So on the whole interval it is at most $e$. Hence

$$\sum_{k=1}^{\lfloor n/2 \rfloor} \left(e^{1-\omega(n)}\frac{e^{pk}}{k}\right)^k \leq \sum_{k=1}^{\lfloor n/2 \rfloor} \left(e^{2-\omega(n)}\right)^k \leq \frac{e^{2-\omega(n)}}{1 - e^{2-\omega(n)}} \to 0.$$

# 3. Polynomials

## 3.1 Schwartz–Zippel Lemma

**Theorem 3.1.1** (Schwartz–Zippel)*. Let $\mathbb{F}$ be an arbitrary field. Let $S$ be a finite subset of $\mathbb{F}$. Suppose that $p(x_1, \ldots, x_m)$ is a polynomial of degree $d$ with coefficients from $\mathbb{F}$. Then the number of $(s_1, \ldots, s_m) \in S^m$ with $p(s_1, \ldots, s_m) = 0$ is at most $d|S|^{m-1}$. In other words, if we choose $s_1, \ldots s_m \in S$ independently and uniformly at random, then the probability that $p(s_1, \ldots, s_m) = 0$ is at most $\frac{d}{|S|}$.*

*Proof.* We prove the claim by induction on $m$. For $m = 1$ the statement claims that a univariate degree $d$ polynomial has at most $d$ zeros, this is well-known. Now suppose that $m > 1$. Let us write $p(x_1, \ldots, x_m)$ in the following form:

$$p(x_1, \ldots, x_m) = \sum_{j=0}^{k} p_j(x_1, \ldots, x_{m-1}) x_m^j,$$

where $k = \deg_{x_m} p$. Note that $\deg p_k(x_1, \ldots, x_{m-1}) = d - k$. Let

$$S_0 = \{(s_1, \ldots, s_{m-1}) \mid s_i \in S, p_k(s_1, \ldots, s_{m-1}) = 0\},$$

and

$$S_1 = \{(s_1, \ldots, s_{m-1}) \mid s_i \in S, p_k(s_1, \ldots, s_{m-1}) \neq 0\},$$

By induction on $m$ we have $|S_0| \leq (d-k)|S|^{m-2}$. If $(s_1, \ldots, s_{m-1}) \in S_1$, then the polynomial

$$p(s_1, \ldots, s_{m-1}, x_m) = \sum_{j=0}^{k} p_j(s_1, \ldots, s_{m-1}) x_m^j,$$

has at most $k$ solutions. Hence the number of solutions of $p(s_1, \ldots, s_m) = 0$ with $s_1, \ldots, s_m \in S$ is at most $|S_0| \cdot |S| + |S_1|k \leq (d-k)|S|^{m-2} \cdot |S| + |S|^{m-1} \cdot k = d|S|^{m-1}$. We are done.

$\square$

## 3.2 Perfect matchings in bipartite graphs

In this section we show how we can use the Schwartz-Zippel lemma to decide whether a bipartite graph contains a perfect matching. Suppose that $G = (A, B, E)$ is a bipartite graph such that $|A| = |B| = n$. For sake of simplicity we assume that the elements of $A$ and $B$ are labelled by the elements of $\{1, 2, \ldots, n\}$. Let us introduce the matrix $R$ of size $n \times n$ as follows: $R_{ij} = x_{ij}$ if $i \in A$ and $j \in B$ are adjacent, and $R_{ij} = 0$ if they are not adjacent. Here $x_{ij}$ is just a variable. Note that if $G$ does not contain a perfect matching, then $\det(S) = 0$. If it contains a perfect matching, say $M$, then nothing cancels the term $(-1)^s \prod_{(i,j) \in M} x_{ij}$ in the expansion of $\det(S)$. In this case $\det(S)$ is a multivariate polynomial of degree $n$. Note that we cannot use Gauss elimination to a matrix containing variables (why?), but we can do the following: we randomly substitute elements of $S$ into $x_{ij}$ and check whether the determinant is non-zero or not. If $\det(R) \neq 0$, then the probability that after the evaluation the result is 0 is at most $\frac{n}{|S|}$. So choose a set $S$ of size $4n$ and do the following algorithm: pick random elements of $S$ and evaluate $\det(R)$. If it is non-zero, then $G$ has a perfect matching. If it is 0, then output that it has no perfect matching. The probability that the algorithm errs, that is, it has a perfect matching, is at most $1/4$. Iterating this process $t$ times the probability that the algorithm errs is at most $1/4^t$.

One detail that might be interesting is that it is worth choosing the set $S$ in a finite field $\mathbb{F}_p$. This way we can save the trouble with counting with fractions or with large numbers. So we choose a prime $p$ bigger than $4n$, and we can even choose $S$ to be the whole $\mathbb{F}_p$.

# 4.  Generating Functions

## 4.1  Two ways of generating functions

The method of generating functions is a powerful tool in enumerative combinatorics. It allows to obtain combinatorial identities by simple algebraic manipulations. The generating function of a sequence $(a_n)$ is simply

$$\sum_{n=0}^{\infty} a_n x^n.$$

One can look at this expression in two different ways. On the one hand, this is just an analytic object, a power series or if you wish, a Taylor-series. On the other hand, one can look at it as a purely algebraic object, the element of the ring

$$\left\{ \sum_{n=0}^{\infty} c_n x^n \mid c_n \in \mathbb{F} \right\},$$

where $\mathbb{F}$ is some field. Both approaches have advantages: when we think of it as an analytic object we can use all our knowledge from analysis. When we consider it as an algebraic object we don't have to check every time the convergence of the series. The element $\sum_{n=0}^{\infty} n! x^n$ is a perfectly acceptable element of the above ring in spite of the fact that it never converges for $|x| > 0$. Honestly, we will do a sloppy, but justifiable thing in this chapter: we consider them as algebraic elements and don't check the convergence, but at the same time we use our analysis knowledge. For instance, we will use that $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$ without elaborating on the question what does it actually mean in an algebraic sense. We will even do differentiation and integration, and this can be formalized in an algebraic way again, but we will not care about building up a formal theory, we simply accept that it works in analysis so it should work here too. One thing that we don't do is that we don't plug any value in a formal power series since we never actually checked the convergence.

In order to use generating functions we need the closed form of some power series. One of the most important is

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n = 1 + x + x^2 + x^3 + \ldots$$

Binomial theorem gives a finite generating function:

$$(1+x)^n = \sum_{k=0}^{n} \binom{n}{k} x^k.$$

There is a similar generating function where the running parameter in the binomial coefficient is the top one:

$$\frac{x^m}{(1-x)^{m+1}} = \sum_{k=0}^{\infty} \binom{k}{m} x^k.$$

(Note that $\binom{k}{m} = 0$ if $0 \le k < m$.) You will prove this result at the recitation. The more power series you learn, the more chance you have for obtaining identities just by algebraic manipulations.

I feel a bit awkward to mention it, but experience shows that even at Msc students I have to say a few words about algebraic manipulations. Algebraic manipulations requires some skill (but not too much knowledge). For instance, one needs to change the order of summations frequently. Please, make sure that you can change the order of the summations in a double sum like

$$\sum_{m=0}^{\infty} \sum_{k=0}^{m} mk^2 x^k.$$

If you are unsure what to do, then draw a 2-dimensional table with running parameters $m$ and $k$, and check that at which field there is a non-zero element, and what happens if you switch to summation for columns instead of rows. Another related advice that if you have a recursion formula for a sequence $(a_n)$ and you have to determine the generating function $\sum a_n x^n$, then it is worth writing out the first few elements separately, because the recursion may act differently there.

Suggested reading: H. Wilf: generatingfunctionology [12]. This book is available online.

### 4.1.1 Exponential generating functions

Besides studying ordinary generating functions sometimes it will be more convenient to work with the so-called exponential generating function of a sequence $(a_n)$:

$$\sum_{n=0}^{\infty} a_n \frac{x^n}{n!}.$$

It requires some skill to decide that in a certain problem we need to use ordinary or exponential generating functions. A strong hint that we need to use exponential generating function is that the sequence $(a_n)$ grows faster then $C^n$ for any $C$, for instance $a_n = n!$. Some recursions also give a hint that it is better to use exponential generating functions.

## 4.2 Enumeration theory and generating functions

The goal of this section is to show some examples of generating functions in enumerative combinatorics. First we will study the so-called Bell-numbers. The Bell-number $B_n$ counts the number of ways to decompose the set $\{1, 2, \ldots, n\}$ into non-empty subsets. For instance, $B_3 = 5$ since the set $\{1, 2, 3\}$ can be decomposed as follows: $\{1, 2, 3\}; \{1, 2\}\{3\}; \{1, 3\}\{2\}; \{2, 3\}\{1\}; \{1\}\{2\}\{3\}$. It is worth defining $B_0$ to be 1.

**Theorem 4.2.1.** *Let $B_n$ denote the number of ways one can decompose the set $\{1, 2, \ldots, n\}$.*

*(a) Then $B_n$ satisfies the recursion*

$$B_n = \sum_{k=0}^{n} \binom{n-1}{k} B_k.$$

*(b) Let us consider the exponential generating function $B(x) = \sum_{n=0}^{\infty} B_n \frac{x^n}{n!}$. Then*

$$B(x) = e^{e^x - 1}$$

*(c) We have*

$$B_n = \frac{1}{e} \sum_{k=1}^{\infty} \frac{k^n}{k!}.$$

*Proof.* (a) Let us choose $k$ elements from $\{1, 2, \ldots, n-1\}$ that are not in the same set of the partition that contains the element $n$. We can do it $\binom{n-1}{k}$ ways, and then we can partition it in $B_k$ ways. (Note that $B_0 = 1$ corresponds to the partition into the single set $\{1, 2, \ldots, n\}$.) Hence

$$B_n = \sum_{k=0}^{n} \binom{n-1}{k} B_k.$$

(b) First we show that $B(x)e^x = B'(x)$ by using the recursion proved in part (a). Note that

$$B(x)e^x = \left(\sum_{n=0}^{\infty} B_n \frac{x^n}{n!}\right) \left(\sum_{n=0}^{\infty} \frac{x^n}{n!}\right) = \sum_{n=0}^{\infty} \left(\sum_{k=0}^{n} \binom{n}{k} B_k\right) \frac{x^n}{n!} = \sum_{n=0}^{\infty} B_{n+1} \frac{x^n}{n!}.$$

On the other hand,

$$B'(x) = \left(\sum_{n=0}^{\infty} B_n \frac{x^n}{n!}\right)' = \sum_{n=1}^{\infty} B_n \frac{x^{n-1}}{(n-1)!} = \sum_{n=0}^{\infty} B_{n+1} \frac{x^n}{n!}.$$

Hence $B(x)e^x = B'(x)$. Next we solve this differential equation: $\frac{B'(x)}{B(x)} = e^x$. Here $\frac{B'(x)}{B(x)} = (\ln B(x))'$, so $\ln B(x) = e^x + c$, thus $B(x) = e^{e^x + c}$. To find the correct value of $c$ observe that $B(0) = 1$, and so $c = -1$. Hence

$$B(x) = e^{e^x - 1}.$$

(c) We have

$$B(x) = \sum_{n=0}^{\infty} B_n \frac{x^n}{n!} = e^{e^x - 1} = \frac{1}{e} e^{e^x} = \frac{1}{e} \sum_{k=0}^{\infty} \frac{e^{kx}}{k!} = \frac{1}{e} \sum_{k=0}^{\infty} \frac{1}{k!} \sum_{n=0}^{\infty} \frac{(kx)^n}{n!} = \sum_{n=0}^{\infty} \frac{x^n}{n!} \left(\frac{1}{e} \sum_{k=0}^{\infty} \frac{k^n}{k!}\right)$$

Hence

$$B_n = \frac{1}{e} \sum_{k=1}^{\infty} \frac{k^n}{k!}.$$

$\square$

**Remark 4.2.2.** I really like Theorem 4.2.1 as it shows a compact story how generating functions work. First we establish a recursion formula from the combinatorial definition. Then we turn it into a generating function. Finally with some algebraic manipulation we deduce some identity that is on the verge of black magic.

Next we study the so-called Stirling-numbers of second kind. This is a refinement of the Bell-numbers.

Let $\left\{ {n \atop k} \right\}$ denote the number of ways to decompose the set $\{1, 2, \ldots, n\}$ into exactly $k$ non-empty sets. For instance, we have $\left\{ {n \atop 1} \right\} = 1$ and $\left\{ {n \atop n} \right\} = 1$ for all $n$, and for instance $\left\{ {3 \atop 2} \right\} = 3$. Clearly,

$$B_n = \sum_{k=1}^{n} \left\{ {n \atop k} \right\}.$$

In order to be able to compute various generating functions of the Stirling-numbers we need some recursion formula for them.

**Proposition 4.2.3.** *We have*

$$\left\{ {n \atop k} \right\} = \left\{ {n-1 \atop k-1} \right\} + k \left\{ {n-1 \atop k} \right\}.$$

*Proof.* Let us consider the partitions of the set $\{1, 2, \ldots n\}$ into $k$ non-empty sets. The number of partitions in which the element $n$ itself determine a set is $\left\{ {n-1 \atop k-1} \right\}$ since we can decompose the remaining $n-1$ elements to $k-1$ sets in $\left\{ {n-1 \atop k-1} \right\}$ ways. If the element $n$ does not determine a set itself then we can decompose the remaining $n-1$ elements into $k$ parts and we can put the element $n$ into any of the $k$ sets. Hence

$$\left\{ {n \atop k} \right\} = \left\{ {n-1 \atop k-1} \right\} + k \left\{ {n-1 \atop k} \right\}.$$

$\square$

**Proposition 4.2.4.** *We have*

$$\sum_{k} \left\{ {n \atop k} \right\} x(x-1)\ldots(x-k+1) = x^n.$$

*First proof.* We prove this identity by induction on $n$ using the recursion formula of the previous proposition:

$$\left\{ {n \atop k} \right\} = \left\{ {n-1 \atop k-1} \right\} + k \left\{ {n-1 \atop k} \right\}.$$

For $n = 1$ the statement is trivial. Then

$$RHS = \sum_{k=1}^{n} \left\{ {n \atop k} \right\} x(x-1)\dots(x-k+1)$$

$$= \sum_{k=1}^{n} \left( \left\{ {n-1 \atop k-1} \right\} + k \left\{ {n-1 \atop k} \right\} \right) x(x-1)\dots(x-k+1)$$

$$= \sum_{k=1}^{n} \left\{ {n-1 \atop k-1} \right\} x(x-1)\dots(x-k+1) + \sum_{k=1}^{n} k \left\{ {n-1 \atop k} \right\} x(x-1)\dots(x-k+1)$$

$$= \sum_{k=1}^{n} \left\{ {n-1 \atop k} \right\} x(x-1)\dots(x-k) + \sum_{k=1}^{n} k \left\{ {n-1 \atop k} \right\} x(x-1)\dots(x-k+1)$$

$$= \sum_{k=1}^{n} \left\{ {n-1 \atop k} \right\} (x(x-1)\dots(x-k) + kx(x-1)\dots(x-k+1))$$

$$= \sum_{k=1}^{n} \left\{ {n-1 \atop k} \right\} x(x-1)\dots(x-k+1)((x-k)+k)$$

$$= x \sum_{k=1}^{n} \left\{ {n-1 \atop k} \right\} x(x-1)\dots(x-k+1)$$

$$= x \cdot x^{n-1}$$

$$= x^n.$$

$\square$

*Second proof.* First let us prove the statement for positive integers $x$. Let us color each element of the set $\{1, 2, \dots, n\}$ with $x$ colors. The number of such coloring is clearly $x^n$. On the hand, we can count this colorings as follows. First we decompose the set into $k$ sets, these will be the color classes. We can color the first set in $x$ ways, the second one in $x - 1$ ways, the third one in $x - 2$ ways... By summing up this for all $k$ we get that

$$\sum_{k} \left\{ {n \atop k} \right\} x(x-1)\dots(x-k+1) = x^n.$$

Since both sides is a polynomial that agree on all positive integers, they must be the same polynomial. $\square$

Once we have set up a recursion we can start proving closed forms for various generating functions.

**Proposition 4.2.5.** *For all $k \geq 1$ we have*

$$\sum_{n \geq 0} \left\{ {n \atop k} \right\} x^n = \frac{x^k}{(1-x)(1-2x)\dots(1-kx)}.$$

*Proof.* Let

$$F_k(x) = \sum_{n \geq 0} \left\{ {n \atop k} \right\} x^n.$$

Using the recursion formula we have

$$F_k(x) = \sum_{n \geq 0} \left\{ {n \atop k} \right\} x^n = \sum_{n \geq 0} \left( \left\{ {n-1 \atop k-1} \right\} + k \left\{ {n-1 \atop k} \right\} \right) x^n = x F_{k-1}(x) + kx F_k(x).$$

Hence

$$F_k(x) = \frac{x}{1-kx} F_{k-1}(x).$$

As we have mentioned $\left\{ {n \atop 1} \right\} = 1$ by definition, so $F_1(x) = \frac{x}{1-x}$. Thus we can conclude that

$$\sum_{n \geq 0} \left\{ {n \atop k} \right\} x^n = \frac{x^k}{(1-x)(1-2x)\dots(1-kx)}.$$

$\square$

In the same vain we can determine the exponential generating functions of the Stirling-numbers of the second kind.

**Proposition 4.2.6.** *For all $k \geq 1$ we have*

$$\sum_{n \geq 0} \left\{ {n \atop k} \right\} \frac{z^n}{n!} = \frac{(e^z - 1)^k}{k!}.$$

*Proof.* Let

$$F_k(z) = \sum_{n \geq 0} \left\{ {n \atop k} \right\} \frac{z^n}{n!}.$$

Then

$$F'_k(z) = \sum_{n \geq 0} \left\{ {n \atop k} \right\} \frac{z^{n-1}}{(n-1)!} = \sum_{n \geq 0} \left( \left\{ {n-1 \atop k-1} \right\} + k \left\{ {n-1 \atop k} \right\} \right) \frac{z^{n-1}}{(n-1)!} =$$

$$= F_{k-1}(z) + k F_k(z).$$

Set $G_k(z) = F_k(z) e^{-kz}$. Then

$$G'_k(z) = (F'_k(z) - k F_k(z)) e^{-kz} = F_{k-1}(z) e^{-kz}.$$

From this it follows by induction that

$$\sum_{n \geq 0} \left\{ {n \atop k} \right\} \frac{z^n}{n!} = \frac{(e^z - 1)^k}{k!}.$$

$\square$

Clearly, if there are Stirling-numbers of the second kind then there must be Stirling-numbers of the first kind. Indeed, the Stirling-numbers of the first kind $\left[ {n \atop k} \right]$ counts the number of permutations of $n$ elements with exactly $k$ cycles in its cycle representation. For instance, $\left[ {4 \atop 2} \right] = 11$ since the following permutations have 2 cycles in its cycle representation: $(1)(234), (1)(243), (2)(134), (2)(143), (3)(124),$ $(3)(142), (4)(123), (4)(132), (12)(34), (13)(24), (14)(23)$. Clearly,

$$\sum_{k=1}^{n} \left[ {n \atop k} \right] = n!.$$

Below we give the analogs of the above statements for the Stirling-number of the first kind without proofs. It's a good exercise for the recitation to prove these claims.

**Proposition 4.2.7.** *For all $k \geq 1$ we have*

$$\left[ {n \atop k} \right] = \left[ {n-1 \atop k-1} \right] + (n-1) \left[ {n-1 \atop k} \right].$$

**Proposition 4.2.8.** *For all $n \geq 1$ we have*

$$\sum_{k=0}^{n} \left[ {n \atop k} \right] x^k = x(x+1) \ldots (x+n-1).$$

**Proposition 4.2.9.** *We have*

$$\sum_{n \geq 0} \left[ {n \atop k} \right] \frac{z^n}{n!} = \frac{1}{k!} \left( \log \frac{1}{1-z} \right)^k.$$

Finally, the following statement connects the Stirling-numbers of the first kind with the Stirling-numbers of the second kind.

**Proposition 4.2.10.** *For all integers $m, n \geq 0$ we have*

$$\sum_{k} \left\{ {n \atop k} \right\} \left[ {k \atop m} \right] (-1)^{n-k} = \begin{cases} 1 & \text{if } m = n, \\ 0 & \text{if } m \neq n. \end{cases}$$

## 4.3 Snake oil method

Generating functions also provide a powerful tool to evaluate certain sums. The so-called snake oil method is very simple, yet handles various sums. Roughly, the idea is the following. Suppose we have some sum that we would like to evaluate, say $\sum_{k=0}^{n} \binom{n}{k}$. Let us call it $A_n$. Then we determine $\sum_{n=0}^{\infty} A_n x^n$: generally this requires a change of summation and some simple algebraic manipulation. Once we have the generating function, in our case $\frac{1}{1-2x}$, we start to determine its coefficients: $\frac{1}{1-2x} = \sum_{n=0}^{\infty} 2^n x^n$. From this we conclude that $A_n = 2^n$. Below you can find several examples for this strategy. Can you fill the gaps in the above argument?

**Proposition 4.3.1.** *We have*

$$\sum_{k=0}^{n} \binom{n+k}{2k} 2^{n-k} = \frac{1}{3}(2 \cdot 4^n + 1).$$

*Proof.* Let

$$A_n = \sum_{k=0}^{n} \binom{n+k}{2k} 2^{n-k}.$$

Then

$$\begin{aligned}
\sum_{n=0}^{\infty} A_n x^n &= \sum_{n=0}^{\infty} \left( \sum_{k=0}^{n} \binom{n+k}{2k} 2^{n-k} \right) x^n \\
&= \sum_{k=0}^{\infty} \frac{1}{2^k} \left( \sum_{n} \binom{n+k}{2k} (2x)^n \right) \\
&= \sum_{k=0}^{\infty} \frac{1}{2^k} \frac{(2x)^k}{(1-2x)^{2k+1}} \\
&= \frac{1}{1-2x} \sum_{k=0}^{\infty} \left( \frac{x}{(1-2x)^2} \right)^k \\
&= \frac{1}{1-2x} \frac{1}{1 - \frac{x}{(1-2x)^2}} \\
&= \frac{1-2x}{1-5x+4x^2} = \frac{1-2x}{(1-x)(1-4x)} \\
&= \frac{2}{3} \frac{1}{1-4x} + \frac{1}{3} \frac{1}{1-x} \\
&= \frac{2}{3} \sum_{n} (4x)^n + \frac{1}{3} \sum_{n} x^n.
\end{aligned}$$

Hence

$$\sum_{k=0}^{n} \binom{n+k}{2k} 2^{n-k} = \frac{1}{3}(2 \cdot 4^n + 1).$$

$\square$

**Proposition 4.3.2.** *We have*

$$\sum_{k} \binom{m}{k}\binom{n+k}{m} = \sum_{k} \binom{m}{k}\binom{n}{k} 2^k.$$

*Proof.* Set

$$A_n = \sum_{k} \binom{m}{k}\binom{n+k}{m},$$

and

$$B_n = \sum_{k} \binom{m}{k}\binom{n}{k} 2^k$$

Then

$$\sum_{n=0}^{\infty} A_n x^n = \sum_{n=0}^{\infty} \left( \sum_{k} \binom{m}{k}\binom{n+k}{m} \right) x^n$$

$$= \sum_{k=0}^{\infty} \binom{m}{k} \left( \sum_{n} \binom{n+k}{m} x^n \right)$$

$$= \sum_{k=0}^{\infty} \binom{m}{k} \frac{x^{m-k}}{(1-x)^{m+1}}$$

$$= \frac{x^m}{(1-x)^{m+1}} \sum_{k} \binom{m}{k} x^{-k}$$

$$= \frac{x^m}{(1-x)^{m+1}} \left( 1 + \frac{1}{x} \right)^m$$

$$= \frac{(1+x)^m}{(1-x)^{m+1}}.$$

On the other hand,

$$\sum_{n=0}^{\infty} B_n x^n = \sum_{n=0}^{\infty} \left( \sum_k \binom{m}{k}\binom{n}{k} 2^k \right) x^n$$

$$= \sum_{k=0}^{\infty} \binom{m}{k} 2^k \left( \sum_n \binom{n}{k} x^n \right)$$

$$= \sum_{k=0}^{\infty} \binom{m}{k} 2^k \frac{x^k}{(1-x)^{k+1}}$$

$$= \frac{1}{1-x} \sum_{k=0}^{\infty} \binom{m}{k} \left( \frac{2x}{1-x} \right)^k$$

$$= \frac{1}{1-x} \left( 1 + \frac{2x}{1-x} \right)^m$$

$$= \frac{(1+x)^m}{(1-x)^{m+1}}.$$

Hence $A_n = B_n$ for all $n$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Proposition 4.3.3.** *We have*

$$\sum_{k=0}^{n} (-1)^{n-k} \binom{2k}{k} \binom{k}{n-k} = 2^n.$$

*Proof.*

$$A_n = \sum_{k=0}^{n} (-1)^{n-k} \binom{2k}{k} \binom{k}{n-k}.$$

Then

$$\sum_{n=0}^{\infty} A_n x^n = \sum_{n=0}^{\infty} \left( \sum_{k=0}^{n} (-1)^{n-k} \binom{2k}{k} \binom{k}{n-k} \right) x^n$$

$$= \sum_{k=0}^{\infty} \binom{2k}{k} x^k \left( \sum_n \binom{k}{n-k} (-x)^{n-k} \right)$$

$$= \sum_{k=0}^{\infty} \binom{2k}{k} x^k (1-x)^k = \sum_{k=0}^{\infty} \binom{2k}{k} (x(1-x))^k$$

$$= \frac{1}{\sqrt{1-4x(1-x)}}$$

$$= \frac{1}{1-2x}$$

$$= \sum_{n=0}^{\infty} 2^n x^n.$$

59

Hence

$$\sum_{k=0}^{n}(-1)^{n-k}\binom{2k}{k}\binom{k}{n-k}=2^n.$$

$\square$

# 5.  Elements of Machine Learning

Machine learning is probably one of the fastest developing part of science with many impressive applications throughout the various aspects of life ranging from voice recognition through spam filters to medical applications. In spite of these fantastic applications the goal of this chapter is not really to teach machine learning, rather an advertisement to your own math knowledge through machine learning. So I picked (rather randomly) elements of machine learning where the underlying mathematics is also interesting.

## 5.1  Classifiers

A large body of machine learning treats the following problem: given examples with labels yes or no (true-false, cat-dog etc) and we try to find a (simple) rule that labels future examples correctly. For instance, a machine is given pictures of cats and dogs with correct labels of cat and dog, and the task is to teach the machine to find a pattern that enables it to label future pictures of cats and dogs properly. There are very many different classifiers: bayesian classifiers, decision tree classifiers, nearest neighbor classifiers, linear and polynomial classifiers, artificial neural networks to mention a few. In the next section we investigate a linear classifier called perceptron algorithm.

### 5.1.1  Linear classifier

Suppose that we would like to buy a vacuum cleaner. We consider three attributes of a vacuum cleaner important: suction power, sound (how noisy it is) and price. We can quantify every vacuum cleaner with three numbers this way, so we can associate a point in $\mathbb{R}^3$ to each vacuum cleaner. Suppose that we label 100 vacuum cleaners whether a consumer would buy it or not. The task of a classifier is to find out the

consumer's taste, that is, label the remaining 10000 vacuum cleaners whether (s)he would buy it or not. Quite often it is possible to find a hyperplane that separates positive examples from negative ones. (Intuitively, humans often consider some (not-quantified) weighted sum that explains their shopping decisions.) So after this long introduction we can formulate it as a mathematical question.

**Problem 5.1.1.** Given $\underline{x}_1, \underline{x}_2, \ldots, \underline{x}_n \in \mathbb{R}^d$ with labels $c(\underline{x}_j) \in \{-1, 1\}$. Suppose that we get the promise that there exists a vector $\underline{c} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ such that

$$\{\underline{x}_j \mid c(\underline{x}_j) = 1\} = \{\underline{x}_j \mid (\underline{c}, \underline{x}_j) > b\} \quad \text{and} \quad \{\underline{x}_j \mid c(\underline{x}_j) = -1\} = \{\underline{x}_j \mid (\underline{c}, \underline{x}_j) < b\}.$$

Problem: find algorithmically (in a fast way!) a $\underline{c}'$ and $b'$ such that

$$\{\underline{x}_j \mid c(\underline{x}_j) = 1\} = \{\underline{x}_j \mid (\underline{c}', \underline{x}_j) > b'\} \quad \text{and} \quad \{\underline{x}_j \mid c(\underline{x}_j) = -1\} = \{\underline{x}_j \mid (\underline{c}', \underline{x}_j) < b'\}.$$

In what follows we will assume that $b = 0$, that is, the separating hyperplane goes through $\underline{0}$, and we would like to find such a hyperplane. The general case can be reduced to this special case by considering the points $(\underline{x}_j, 1) \in \mathbb{R}^{d+1}$, and the last coordinate of $\underline{c} \in \mathbb{R}^{d+1}$ will correspond to $-b$.

Here is the solution of Rosenblatt called perceptron.

**Algorithm 5.1.2.** (Perceptron algorithm) Let $\underline{w} = \underline{0}$.

1. If there exists an $\underline{x}_j$ such that $\text{sign}(\underline{w}, \underline{x}_j) \neq c(\underline{x}_j)$, then let $\underline{w} = \underline{w} + c(\underline{x}_j)\underline{x}_j$.

2. If there exists no $\underline{x}_j$ such that $\text{sign}(\underline{w}, \underline{x}_j) \neq c(\underline{x}_j)$ then output $\underline{c}' = \underline{w}$. Otherwise go to step 1.

It is clear that if the perceptron algorithm halts, then it outputs a separating hyperplane. The question is why it halts and how many steps.

**Theorem 5.1.3.** *Given $\underline{x}_1, \underline{x}_2, \ldots, \underline{x}_n \in \mathbb{R}^d$ with labels $c(\underline{x}_j) \in \{-1, 1\}$. Suppose that we get the promise that there exists a vector $\underline{c} \in \mathbb{R}^d$ such that*

$$\{\underline{x}_j \mid c(\underline{x}_j) = 1\} = \{\underline{x}_j \mid (\underline{c}, \underline{x}_j) > 0\} \quad \text{and} \quad \{\underline{x}_j \mid c(\underline{x}_j) = -1\} = \{\underline{x}_j \mid (\underline{c}, \underline{x}_j) < 0\}.$$

*Let $R = \max_{j=1,\ldots,n} ||x_j||$ and $\gamma = \min_{j=1,\ldots,n} |(\underline{x}_j, \underline{c})|$. Then the perceptron algorithm halts in $\frac{R^2}{|\underline{c}|^2 \gamma^2}$ steps.*

*Proof.* Let $\underline{w}_0 = \underline{0}$, and for a general $k$ if there exists an $\underline{x}_j$ such that $\text{sign}(\underline{w}_k, \underline{x}_j) \neq c(\underline{x}_j)$, then let $\underline{w}_{k+1} = \underline{w}_k + c(\underline{x}_j)\underline{x}_j$. Observe that

$$(\underline{w}_{k+1}, \underline{c}) = (\underline{w}_k, \underline{c}) + c(\underline{x}_j)(\underline{x}_j, \underline{c}) \geq (\underline{w}_k, \underline{c}) + \gamma$$

since $\text{sign}(\underline{x}_j, \underline{c}) = c(\underline{x}_j)$. Secondly,

$$(\underline{w}_{k+1}, \underline{w}_{k+1}) = (\underline{w}_k, \underline{w}_k) + 2c(x_j)(\underline{w}_k, \underline{x}_j) + (\underline{x}_j, \underline{x}_j) < (\underline{w}_k, \underline{w}_k) + (\underline{x}_j, \underline{x}_j) \leq (\underline{w}_k, \underline{w}_k) + R^2$$

since $\text{sign}(\underline{w}_k, \underline{x}_j) \neq c(\underline{x}_j)$ means that $2c(x_j)(\underline{w}_k, \underline{x}_j) < 0$. So after $M$ steps we get that

$$(\underline{w}_M, \underline{c}) > M\gamma \quad \text{and} \quad (\underline{w}_M, \underline{w}_M) < MR^2.$$

Now let us use the Cauchy-Schwarz inequality:

$$MR^2(\underline{c}, \underline{c}) > (\underline{w}_M, \underline{w}_M)(\underline{c}, \underline{c}) \geq (\underline{w}_M, \underline{c})^2 \geq M^2\gamma^2.$$

Hence $M < \frac{R^2|\underline{c}|^2}{\gamma^2}$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

As we have seen the above algorithm directly turns into a machine learning problem. Suppose we want to find out the taste of a consumer on vacuum cleaners in order to sell him the most expensive one. We show him a few vacuum cleaners and ask him/her to label it whether he would buy it or not. Then we choose a hyperplane that separates the positive examples from negative examples. Then we use this hyperplane for our prediction to unlabeled vacuum cleaners (and then we try to sell the most expensive one that is predicted to be vendible). Note that we never said that it is the correct hyperplane, we only said that we will use this hyperplane for our prediction.

Let's phrase it in a machine learning language. We had a training set $S$ (later it will be convenient to think of it as a sequence of examples) that were labeled and our task was to label future examples. This was a **classification** problem, we trained a **classifier**. Another type of machine learning problems is **regression**, where we have to predict a value and not a label. This was a **supervised learning** since we got the labels. This was also **batch learning protocol** which means that we got a lot of examples and we needed to act only after studying them, the opposite is when we have to act **online** so the learning and decision making happens at the same time like in a case of a stockbroker. This was a **passive learner** since we had no chance to select what examples are to be labelled, an **active learner** can select which examples should be labelled.

What is more interesting that we assumed that a half-space gives the correct labeling and we were looking for our prediction already in this form. These are two different things: we can search for the solution in a form of hyperplane even if nobody guarantees that the ultimately correct labeling is in this form. In machine learning

language we had a hypothesis class $\mathcal{H}$, the set of half-spaces in our example. So our output of the learning was a $h \in \mathcal{H}$, a **prediction rule** (also called **predictor, hypothesis** or **classifier**). One might wonder why we made this assumption. We will see later that it is absolutely necessary for learning that we have some preliminary assumptions on the problem. A final observation is that even though we were promised to have an optimal hyperplane it was still an **algorithmic challenge** to find such a hyperplane, and of course nothing guarantee that it is the optimal hyperplane and future examples will be labelled properly by the prediction rule we have found.

## 5.1.2   Formal model for the statistical learning framework

In this section we try to formalize what we mean by statistical learning.

**Domain set:** An arbitrary set $\mathcal{X}$. This is the set of objectes that we wish to label. Usually, these domain points will be represented by a vector of **features**. We will also call the elements of $\mathcal{X}$ instances, and $\mathcal{X}$ the instance space.

**Label set:** It can be $\{yes, no\}$ or $\{dog, cat, bird\}$ so any set. We will denote by $\mathcal{Y}$ the set of possible labels. Actually, we can use this framework for regression problems too when we have to predict a value instead of a label, so $\mathcal{Y}$ can $\mathcal{R}$ too.

**Training data:** A sequence $S = ((x_1, y_1), \ldots, (x_m, y_m))$ is a finite sequence of pairs in $\mathcal{X} \times \mathcal{Y}$. So this is a sequence of labeled domain points. This is the input of the learner. These are also called training examples or training set (although we will treat them as a sequence).

**The learner's output:** This is a function $h : \mathcal{X} \to \mathcal{Y}$. It is called **prediction rule, predictor, hypothesis, classifier**. If we want to emphasize that this function was produced by an algorithm $A$ based on the training set $S$, then we denote it by $A(S)$. We will often assume that there is a hypothesis class $\mathcal{H}$ from which we choose the function $h$.

**Measuring success:** Let us introduce a **loss function** $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \to \mathbb{R}_{\geq 0}$ that measures how successful the hypothesis $h$ on an element $(x, y) \in \mathcal{X} \times \mathcal{Y}$. For instance, if $\mathcal{Y}$ is a set of labels, then we may say that

$$\ell_{0-1}(h, (x, y)) = \begin{cases} 0 & \text{if} h(x) = y \\ 1 & \text{if} h(x) \neq y \end{cases}$$

In the continuous case (regression problem), when the prediction rule outputs a number we might use the loss function $\ell((h, (x, y)) = (h(x) - y)^2$.

Next we define the **risk function**. Here we assume that there is a probability distribution $\mathcal{D}$ on $\mathcal{X} \times \mathcal{Y}$. For instance, this might be the uniform distribution on all vacuum cleaners (with a hypothetical label on it). Then

$$L_{\mathcal{D}}(h) := \mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(h, (x, y)).$$

Note that $\mathcal{D}$ is not known. Also note that it can occur that $\mathbb{P}_{\mathcal{D}}(x, y_1) > 0$ and $\mathbb{P}_{\mathcal{D}}(x, y_2) > 0$ with different $y_1$ and $y_2$. For instance, if there are two vacuum cleaners with the same parameters (suction power, sound and price), but with different colors and our costumer would only buy the blue one, then it results an $(x, y)$ for which $1 > \mathbb{P}_{\mathcal{D}}(x, y) > 0$. We also have an **empirical risk**:

$$L_S(h) := \frac{1}{m} \sum_{i=1}^{m} \ell(h, (x_i, y_i)).$$

**Empirical risk minimization:** In many cases we will assume that there is an (efficient) algorithm that outputs an $h \in \mathcal{H}$ minimizing $L_S(h)$. Note that there can be many such $h \in \mathcal{H}$. Our hope is that if $m$ is large (the size of the training data), then $L_{\mathcal{D}}(h)$ will not much larger than $L_S(h)$. Formalizing this idea will be the content of the next section.

**Realizability assumption:** In the above example we assumed that there exists a hyperplane that separates all instances well, not just the training set, that is, there exists an $h^*$ such that $L_{\mathcal{D}}(h^*) = 0$. This is called realizability assumption. It often simplifies mathematical examination, but not a natural assumption in real world applications.

## 5.2   PAC-learning

In this section we elaborate on the problem when we can expect that a prediction rule output by the empirical risk minimization is indeed a good predictor. So what can go wrong? The first thing that we might think of that the training set $S$ was not very similar to the distribution $\mathcal{D}$ or did not catch some important characteristic of

$\mathcal{D}$. This can always happen even if with not too high probability. So we can never be absolutely sure that our prediction rule will work well on $\mathcal{D}$ too. Another important observation is that the role of $\mathcal{H}$ is crucial. If it contains too few functions, then there is no chance for a good learner. But it is also a problem if it contains too many functions: there will be minimizers that works well on the training set $S$, but badly on all examples. For instance, if $\mathcal{H}$ contains all functions, then the function $h$ with

$$h(x) = \begin{cases} y & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases}$$

is perfect on $S$, but unlikely to be good on all instances. The theory of PAC-learnability addresses these issues. PAC stands for Probably Approximately Correct. Here the word probably refers to the phenomenon that $S$ might have been uncharacteristic, approximately correct refers to the phenomenon that we cannot expect a perfect predictor even if $S$ was a good sample. The formal definition is as follows.

**Definition 5.2.1.** (Agnostic PAC-learnability for general loss function) A hypothesis class $\mathcal{H}$ is agnostic PAC-learnable with respect to $\mathcal{X} \times \mathcal{Y}$ and a loss function $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \to \mathbb{R}_{\geq 0}$ if there exists a function $m_{\mathcal{H}} : (0,1)^2 \to \mathbb{Z}_{\geq 0}$ and a learning algorithm with the following property: for every $\varepsilon, \delta \in (0,1)$ and for every distribution $\mathcal{D}$ on $\mathcal{X} \times \mathcal{Y}$, when running the algorithm on $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ examples generated by $\mathcal{D}$, the algorithm outputs an $h \in \mathcal{H}$ such that with probability at least $1 - \delta$ (over the choice of the $m$ training examples) we have

$$L_{\mathcal{D}}(h) \leq \inf_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \varepsilon.$$

## 5.2.1 VC-dimension

In this section we introduce the concept of Vapnik-Chervonenkis dimension or as it is more commonly called, the VC-dimension.

**Definition 5.2.2.** An $S = (X, R)$ is called a range space, where $X$ is a finite or infinite set, and $R$ is a finite or infinite family of subsets of $X$. The elements of $X$ are called points, the elements of $R$ are called ranges.

If $A \subseteq X$, then $P_R(A) = \{r \cap A | \ r \in R\}$ is the projection of $R$ on $A$.

We say that $A$ is shattered if $P_R(A)$ contains all subsets of $A$. The VC-dimension of $S$ –denoted by $VC(S)$– is the maximum cardinality of a shattered subset.

**Example 5.2.3.** Let $X = \mathbb{R}^2$, and let $R$ contain all half planes. Then three points in general position can be shattered, but no matter how we take 4 points we cannot shatter it. Hence $VC(\mathbb{R}^2, \text{halfplanes}) = 3$.

**Theorem 5.2.4** (Sauer). *Let $(X, R)$ be a range space of VC-dimension $d$ with $|X| = n$. Then $|R| \leq \sum_{k=0}^{d} \binom{n}{k}$.*

*Proof.* Let us introduce the notation $g(n, d) = \sum_{k=0}^{d} \binom{n}{k}$. Note that

$$g(n, d) = g(n - 1, d) + g(n - 1, d - 1).$$

We prove the claim by induction on $n$. For $n = 1$ the claim is trivial. So assume that the claim already holds till $n - 1$. So let $S = (X, R)$ be a range space of VC-dimension $d$ with $|X| = n$. Let $x \in X$ arbitrary, and let us consider the following two range spaces: $S - x = (X - \{x\}, R - x)$ and $S/x = (X - \{x\}, R/x)$, where

$$R - x = \{r \setminus \{x\} \mid r \in R\} \quad \text{and} \quad R/x = \{r \in R \mid x \notin r, \; r \cup \{x\} \in R\}.$$

Observe that the VC-dimension of $S - x$ is at most $d$ while the VC-dimension of $S/x$ is at most $d - 1$ (why?). Hence by induction

$$|R| = |R - x| + |R/x| \leq g(n - 1, d) + g(n - 1, d - 1) = g(n, d).$$

$\square$

**Definition 5.2.5.** Let $(X, R)$ be a range space, and let $A \subseteq X$ be finite. For $0 \leq \varepsilon \leq 1$ a subset $B \subseteq A$ is an $\varepsilon$-sample for $A$ if for any range $r \in R$ the inequality

$$\left| \frac{|A \cap r|}{|A|} - \frac{|B \cap r|}{|B|} \right| \leq \varepsilon.$$

A subset $N \subseteq A$ is an $\varepsilon$-net for $A$ if any range $r \in R$ satisfying $|A \cap r| > \varepsilon|A|$ contains at least one point of $N$.

Note that every $\varepsilon$-net is automatically an $\varepsilon$-sample, but the converse is not true.

**Theorem 5.2.6** (Vapnik and Chervonenkis). *There exists a universal positive constant $c$ with the following properties. Let $0 < \varepsilon, \delta < 1$. Let $(X, R)$ be a range space of VC-dimension $d$, and let $A$ be an arbitrary subset of $X$. Let*

$$s \geq \min\left( |A|, \frac{c}{\varepsilon^2} \left( d \ln \frac{d}{\varepsilon} + \ln \frac{1}{\delta} \right) \right).$$

*Then a random subset $B$ of cardinality $s$ of $A$ is an $\varepsilon$-sample with probability at least $1 - \delta$.*

**Theorem 5.2.7** (Haussler and Welzl)**.** *Let $0 < \varepsilon, \delta < 1$. Let $(X, R)$ be a range space of VC-dimension d, and let $A$ be an arbitrary subset of $X$. Let*

$$m \geq \max\left(\frac{4}{\varepsilon} \ln \frac{4}{\delta}, \frac{8d}{\varepsilon} \ln \frac{8d}{\varepsilon}\right).$$

*Then a random subset $N$ of cardinality $m$ of $A$ is an $\varepsilon$-net with probability at least $1 - \delta$.*

We only prove Theorem 5.2.7.

*Proof of Theorem 5.2.7.* Let $N = \{x_1, \ldots, x_m\}$ be a random multi-subset of $A$ obtained by $m$ independent random draws from $A$. Let $E_1$ be the following event:

$$E_1 = \{\exists r \in R \mid |r \cap A| \geq \varepsilon n, r \cap N = \emptyset\}.$$

We need to prove that $\mathbb{P}(E_1) \leq \delta$. Let $T = \{y_1, \ldots, y_m\}$ be another random multi-subset of $A$ obtained by $m$ independent random draws from $A$. Let $E_2$ be the following event:

$$E_1 = \{\exists r \in R \mid |r \cap A| \geq \varepsilon n, r \cap N = \emptyset, |r \cap T| \geq \frac{\varepsilon m}{2}\}.$$

Here $|r \cap T| = |\{i \mid y_i \in r\}|$. We define $|r \cap N|$ and $|r \cap (N \cup T)|$ similarly. The proof of Theorem 5.2.7 relies on the following two lemmas together with a small computation.

**Lemma 5.2.8.** *We have $\mathbb{P}(E_2) \geq \mathbb{P}(E_1)$.*

**Lemma 5.2.9.** *We have $\mathbb{P}(E_2) \leq g(2m, d)2^{\varepsilon m/2}$.*

*Proof of Lemma 5.2.8.* Since $\mathbb{P}(E_2)/\mathbb{P}(E_1) \geq \mathbb{P}(E_2 \cap E_1)/\mathbb{P}(E_1) = \mathbb{P}(E_2|E_1)$ it is enough to prove that $\mathbb{P}(E_2|E_1) \geq \frac{1}{2}$. Suppose that for $N = \{x_1, \ldots, x_m\}$ the event $E_1$ satisfies. Let us fix an $r$ such that $|r \cap A| \geq \varepsilon n$ and $r \cap N = \emptyset$. It is enough to show that already for this fixed $r$ we have $\mathbb{P}(|r \cap T| \geq \frac{\varepsilon m}{2}) \geq \frac{1}{2}$. Let $p = \frac{|r \cap A|}{|A|} \geq \varepsilon$ and $X = |r \cap T|$. Then $\mathbb{E}X = pm$ and $\mathrm{Var}(X) = p(1-p)m < pm$. Furthermore,

$$\mathbb{P}(X < \frac{\varepsilon m}{2}) = \mathbb{P}(\mathbb{E}X - X > (p - \varepsilon/2)m) \leq \mathbb{P}(|\mathbb{E}X - X| > (p - \varepsilon/2)m)$$

$$\leq \frac{\mathrm{Var}(X)}{((p - \varepsilon/2)m)^2} \leq \frac{pm}{(pm/2)^2} = \frac{4}{pm} \leq \frac{1}{2}.$$

Hence $\mathbb{P}(|r \cap T| \geq \frac{\varepsilon m}{2}) \geq \frac{1}{2}$ thereby implying that $\mathbb{P}(E_2|E_1) \geq \frac{1}{2}$. $\qquad\square$

*Proof of Lemma 5.2.9.* We can choose $N$ and $T$ by first choosing $N \cup T$, and then choosing $N$ and $T$ from $N \cup T$. We have

$$\mathbb{P}(E_2) = \sum_{S=N\cup T} \mathbb{P}(S)\mathbb{P}(E_2|S).$$

So it is enough to show that $\mathbb{P}(E_2|S) \leq g(2m,d)2^{\varepsilon m/2}$ for every $S = N \cup T$. For a fixed $S = N \cup T$ let $E_r = \{r \cap N = \emptyset, |t \cap T| \geq \varepsilon m/2\}$. If $r \cap S = r' \cap S$, then $E_r = E_{r'}$. Since the VC-dimension is $d$, we can only have $g(2m,d)$ different $r \cap S$ sets by Theorem 5.2.4. Now fix an $r \in R$ and suppose that $|r \cap S| = s \geq \varepsilon m/2$, then

$$\mathbb{P}(r \cap N = \emptyset \mid N \cup T = S) = \frac{(2m-s)(2m-s+1)\ldots(m-s+1)}{2m(2m-1)\ldots m}$$

$$= \frac{\frac{(2m-s)!}{(m-s)!}}{\frac{(2m)!}{m!}} = \frac{\frac{m!}{(m-s)!}}{\frac{(2m)!}{(2m-s)!}}$$

$$= \frac{m(m-1)\ldots(m-s+1)}{2m(2m-1)\ldots(2m-s+1)} \leq 2^{-s} \leq 2^{-\varepsilon m/2}.$$

Hence $\mathbb{P}(E_2|S) \leq g(2m,d)2^{-\varepsilon m/2}$. $\qquad\square$

By putting together the two lemmas we get that $\mathbb{P}(E_1) \leq 2g(2m,d)2^{-\varepsilon m/2}$. It is enough to show that if $m$ satisfies the condition of the theorem, then $2g(2m,d)2^{-\varepsilon m/2} \leq \delta$. Note that

$$g(2m,d) = \sum_{k=0}^{d} \binom{2m}{k} \leq 2(2m)^d$$

even if $d = 0$ or $1$ and for larger $d$ even tighter inequalities are true. So

$$2g(2m,d)2^{-\varepsilon m/2} \leq 4(2m)^d 2^{-\varepsilon m/2}.$$

So it is enough to prove the inequality $4(2m)^d 2^{-\varepsilon m/2} \leq \delta$. This is equivalent with

$$\frac{2}{\varepsilon}\ln\left(\frac{4}{\delta}\right) + \frac{2d}{\varepsilon}\ln(2m) \leq m.$$

Here $\frac{2}{\varepsilon}\ln\left(\frac{4}{\delta}\right) \leq \frac{m}{2}$ so it is enough to show that $\frac{2d}{\varepsilon}\ln(2m) \leq \frac{m}{2}$. The derivative of the left hand side with respect to $m$ is $\frac{2d}{\varepsilon m} \leq \frac{1}{2}$ by the condition on $m$, so it is enough to check the statement for $m_0 = \frac{8d}{\varepsilon}\ln\left(\frac{8d}{\varepsilon}\right)$. Let $x = \frac{8d}{\varepsilon} \geq 8$. Then $2\ln x \leq x$, and so $\frac{x}{2}\ln(x\ln x) \leq x\ln x$, that is, $\frac{4d}{\varepsilon}\ln(2m_0) \leq m_0$. This completes the proof of the theorem.

$\qquad\square$

## 5.3 Clustering

Clustering is an important problem in the area of unsupervised machine learning. In this section we study the so-called $k$-mean problem.

The mathematical formulation of the problem is quite simple. Let $\underline{x}_1, \ldots, \underline{x}_n \in \mathbb{R}^d$. In the $k$-mean problem the goal is to find centers $\mathcal{C} = \{\underline{c}_1, \ldots, \underline{c}_k\}$ that (approximately) minimizes the sum

$$F(\mathcal{C}) := \sum_{i=1}^n \min_{j=1}^k ||\underline{x}_i - \underline{c}_j||^2.$$

Once we have such centers we can cluster the points $\underline{x}_1, \ldots, \underline{x}_n$ as follows:

$$C_r = \{\underline{x}_i \mid \min_{j=1}^k ||\underline{x}_i - \underline{c}_j|| = ||\underline{x}_r - \underline{c}_j||\}$$

for $r = 1, \ldots, k$. In words, we put the point $\underline{x}_i$ into the cluster $C_r$ if the closest point among the centers $\underline{c}_1, \ldots, \underline{c}_k$ is $\underline{c}_r$. (In case of ties, we can randomly put the point of $\underline{x}_i$ to one of these clusters.)

Probably the most well-known algorithm for the $k$-means problem is Lloyd's algorithm.

**Algorithm 5.3.1** (Lloyd)**.** Initialize centers $\underline{p}_1, \ldots, \underline{p}_k$ by choosing them randomly from $\underline{x}_1, \ldots, \underline{x}_n$. Then iterate the following steps.

1. Form clusters $C_1, \ldots, C_k$ by putting point $\underline{x}_i$ into the cluster $C_r$ if the closest point among the centers $\underline{p}_1, \ldots, \underline{p}_k$ is $\underline{p}_r$.

2. Having clusters $C_1, \ldots, C_k$ let $\underline{p}_r$ be the center of gravity of the points in $C_r$, that is,

$$\underline{p}_r = \frac{1}{|C_r|} \sum_{\underline{x}_i \in C_r} \underline{x}_i.$$

Stop when $\underline{p}_1, \ldots, \underline{p}_k$ does not change and output them as $\underline{c}_1, \ldots, \underline{c}_k$.

Some words about the stopping rule. In every step $F(\mathcal{C})$ strictly decreases, apart from the case when a point has the same distance from two different centers. Since the number of possible configurations is finite, this means that the algorithm eventually stops. Unfortunately, there are cases when the algorithm converge very slowly. (In fact, due to numerical inaccuracies it may even occur that the algorithm does not converge.) So it might be useful to invent alternative stopping rules like stop when $F(\mathcal{C})$ only changes less than some preset $\varepsilon$.

# Bibliography

[1] M. Ajtai, V. Chvátal, M. M. Newborn, and E. Szemerédi, *Crossing-free subgraphs*, North-Holland Mathematics Studies, **60** (1982), pp. 9–12.

[2] N. Alon and J. H. Spencer, *The probabilistic method*, John Wiley & Sons, 2004.

[3] A. E. Brouwer and W. H. Haemers, *Spectra of graphs*, Springer Science & Business Media, 2011.

[4] G. Elekes, *On the number of sums and products*, Acta Arithmetica, **81** (1997), pp. 365–367.

[5] P. Erdős, *Graph theory and probability*, Canad. J. Math, **11** (1959), pp. 34–38.

[6] P. Erdős, *On a problem of graph theory*, Math. Gaz., **47** (1963), pp. 220–223.

[7] P. Erdős and A. Rényi, *On the evolution of random graphs*, Publ. Math. Inst. Hung. Acad. Sci, **5** (1960), p. 43.

[8] J. Matousek, *Thirty-three miniatures: Mathematical and Algorithmic applications of Linear Algebra*, vol. 53, American Mathematical Soc., 2010.

[9] J. Solymosi, *Bounding multiplicative energy by the sumset*, Advances in mathematics, **222** (2009), pp. 402–408.

[10] R. Stanley, *Topics in algebraic combinatorics*, Course notes for Mathematics, **192** (2000).

[11] E. Szemerédi and W. T. Trotter, *Extremal problems in discrete geometry*, Combinatorica, **3** (1983), pp. 381–392.

[12] H. Wilf, *Generatingfunctionology,(1990)*, ISBN: 0-12-751956-4.