

WARWICK MATHEMATICS EXCHANGE

MA3J2

COMBINATORICS II

2024, June 13th

Desync, aka The Big Ree

Contents

Table of Contents	i
1 Projective Planes and Latin Squares	1
1.1 Projective Planes	1
1.2 Finite Projective Planes	3
1.3 Latin Squares	4
2 Error-Correcting Codes	6
2.1 Introduction	6
2.2 Block Codes	8
2.3 Hamming Codes	9
2.4 Shannon's Theorem	13
3 Discrete Geometry	14
3.1 Separation	16
3.2 Extrema	18
3.3 Polyhedra and Polytopes	20
3.4 Polars	21
3.5 Radon's Lemma and Helly's Theorem	25
4 Partially Ordered Sets and Set Systems	28
4.1 Dilworth's Theorem	30
4.2 Covering by Chains	31
4.3 VC Dimension and the Sauer-Shelah Lemma	33
5 Graph Colouring	35
5.1 The Chromatic Polynomial	38
6 Matroids	43
6.1 Rado's Theorem	44
7 Random Graphs	45
7.1 Chromatic Numbers	47
7.2 Connectedness	48

Introduction

ree

Disclaimer: I make *absolutely no guarantee* that this document is complete nor without error. In particular, any content covered exclusively in lectures (if any) will not be recorded here. This document was written during the 2023 academic year, so any changes in the course since then may not be accurately reflected.

Notes on formatting

New terminology will be introduced in *italics* when used for the first time. Named theorems will also be introduced in *italics*. Important points will be **bold**. Common mistakes will be underlined. The latter two classifications are under my interpretation. YMMV.

Content not taught in the course will be outlined in the margins like this. Anything outlined like this is not examinable, but has been included as it may be helpful to know alternative methods to solve problems.

The table of contents above, and any inline references are all hyperlinked for your convenience.

History

First Edition: 2024-05-28*

Current Edition: 2024-06-13

Authors

This document was written by R.J. Kit L., a maths student. I am not otherwise affiliated with the university, and cannot help you with related matters.

Please send me a PM on Discord @Desync#6290, a message in the WMX server, or an email to Warwick.Mathematics.Exchange@gmail.com for any corrections. (If this document somehow manages to persist for more than a few years, these contact details might be out of date, depending on the maintainers. Please check the most recently updated version you can find.)

If you found this guide helpful and want to support me, you can [buy me a coffee!](#)

(Direct link for if hyperlinks are not supported on your device/reader: ko-fi.com/desync.)

*Storing dates in big-endian format is clearly the superior option, as sorting dates lexicographically will also sort dates chronologically, which is a property that little and middle-endian date formats do not share. See ISO-8601 for more details. This footnote was made by the computer science gang.

1 Projective Planes and Latin Squares

1.1 Projective Planes

There is a deep connection between algebra and geometry, but once we move beyond linear algebra, there is a certain inconvenience in the vector spaces in which we do geometry, which stems primarily from an asymmetry between points and lines. Any two points are incident to a single line, but it is not true that any two lines are incident to a single point: they could be parallel, or in higher dimensions, skew.

To resolve this imbalance, we add a point “at infinity” which some lines may be incident to in *projective geometry*.

Consider the vector space K^3 over a field K . Removing the origin, we define an equivalence relation on $K^3 \setminus \{\mathbf{0}_K\}$ by $(a,b,c) \sim (x,y,z)$ if there exists $0 \neq \lambda \in K$ such that $(a,b,c) = \lambda(x,y,z)$. That is, vectors are equivalent up to scaling.

The *projective plane over K* , denoted $K\mathbb{P}^2$ is then the set of equivalence classes of non-zero vectors in K^3 . Equivalently, the points of $K\mathbb{P}^2$ may be viewed as the lines through the origin in K^3 .

Example. Consider the case $K = \mathbb{R}$, giving the *real projective plane* \mathbb{RP}^2 . Each line through the origin in \mathbb{R}^3 intersects the unit sphere at two antipodal points, so we can also view the set of points of \mathbb{RP}^2 as the surface of the sphere with antipodal points identified.

Intuitively, a line in \mathbb{RP}^2 is then just a great circle on the sphere, also with antipodal points identified. Such a great circle also be viewed as the intersection of a plane through the origin with the sphere, which can be characterised by the normal vector $(\lambda, \mu, \nu) \neq \mathbf{0}_K$. The great circle is then the (equivalence class of the) set of points (x,y,z) satisfying

$$\lambda x + \mu y + \nu z = 0$$

where λ, μ , and ν are elements of K that are not all zero. △

We will be studying the discrete analogue of these spaces in which our projective planes have only finitely many points.

If a field K is finite with q elements, then $K^3 \setminus \{\mathbf{0}_K\}$ has $q^3 - 1$ elements. Each equivalence class has $q - 1$ elements, as there are $q - 1$ non-zero elements of K to scale by. So, there are

$$\frac{q^3 - 1}{q - 1} = q^2 + q + 1$$

elements, or *points*, in $K\mathbb{P}^2$.

A *line* in $K\mathbb{P}^2$ is then the set of points (x,y,z) satisfying

$$\lambda x + \mu y + \nu z = 0$$

where λ, μ , and ν are elements of K that are not all zero. That is, a line is the set of points orthogonal to a point (λ, μ, ν) .

Note that this is well-defined due to the bilinearity of the dot product.

Example. Consider the case $K = \mathbb{Z}_2$, the field of two elements $\{0,1\}$ in which $1 + 1 = 0$. In this case, the equivalence classes in $\mathbb{Z}_2\mathbb{P}^2$ are singletons, so the points in the projective planes are given by the seven non-zero vectors

$$(0,0,1), (0,1,0), (0,1,1), (1,0,0), (1,0,1), (1,1,0), (1,1,1)$$

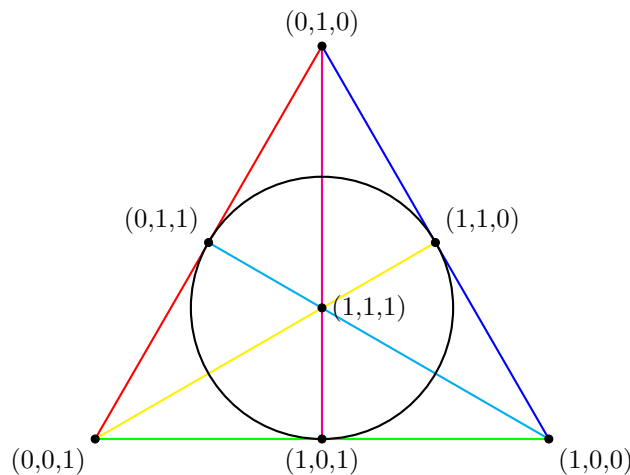
Each line in $\mathbb{Z}_2\mathbb{P}^2$ contains 3 points: for instance, the line represented by $(1,0,0)$ (i.e. $1x + 0y + 0z = 0$) consists of the three points

$$(0,0,1), \quad (0,1,0), \quad (0,1,1)$$

and the line represented by $(1,1,1)$ (i.e. $1x + 1y + 1z = 0$) consists of the three points

$$(0,1,1), \quad (1,0,1), \quad (1,1,0)$$

This projective plane is also called the *Fano plane*, and its points are often drawn arranged in a triangle:



△

As we would expect, any two distinct points determine a unique line connecting them.

Lemma 1.1 (Points in the Projective Plane). *Given any two distinct points in $K\mathbb{P}^2$, there is exactly one line incident to both of them.*

Proof. Let the points be (represented by) (a,b,c) and (x,y,z) and form the cross product

$$\begin{aligned} (\lambda, \mu, \nu) &:= (a,b,c) \times (x,y,z) \\ &= (bz - cy, cx - az, ay - bx) \end{aligned}$$

The cross product is also zero if and only if the two starting vectors are linearly dependent, but by assumption, (a,b,c) and (x,y,z) are representatives of distinct points and are therefore linearly independent, so (λ, μ, ν) is non-zero and defines a line.

By construction, the cross product is orthogonal to (a,b,c) and (x,y,z) (or otherwise, this can be checked by hand), so both the points are incident to the line defined by (λ, μ, ν) .

For uniqueness, suppose (λ', μ', ν') defines a line incident to (a,b,c) and (x,y,z) :

$$\lambda'a + \mu'b + \nu'c = 0 \tag{1}$$

$$\lambda'x + \mu'y + \nu'z = 0 \tag{2}$$

From the above, we have that $bz - cy$, $cx - az$, and $ay - bx$ are not all zero, so without loss of generality, suppose that $\phi := bz - cy \neq 0$. Then, multiplying (1) by z and (2) by c , and subtracting (2) from (1), we have:

$$\begin{aligned} c(\lambda'x + \mu'y + \nu'z) - z(\lambda'a + \mu'b + \nu'c) &= 0 - 0 \\ \lambda'cx + \mu'cy - \lambda'az - \mu'bz &= 0 \end{aligned}$$

$$\begin{aligned}
\lambda'(cx - az) + \mu'(cy - bz) &= 0 \\
\lambda'(cx - az) &= \mu'(bz - cy) \\
\lambda'(cx - az) &= \mu'\phi
\end{aligned}$$

and similarly,

$$\lambda'(ay - bx) = \nu'\phi$$

so

$$\begin{aligned}
\lambda' &= (\phi^{-1}\lambda')(bz - cy) \\
\mu' &= (\phi^{-1}\mu')(cx - az) \\
\nu' &= (\phi^{-1}\nu')(ay - bx)
\end{aligned}$$

so $(\lambda', \mu', \nu') \sim (\lambda, \mu, \nu)$, and the line is unique. ■

The preceding proof shows that if (a, b, c) and (x, y, z) are non-equivalent elements of $K^3 \setminus \{0\}$, then there is a unique equivalence class of elements $(\lambda, \mu, \nu) \in K^3 \setminus \{0\}$ satisfying

$$\lambda a + \mu b + \nu c = \lambda x + \mu y + \nu z = 0$$

However, we may also interpret (a, b, c) and (x, y, z) as (equivalence classes of) lines and $[(\lambda, \mu, \nu)]$ as a point, so this also shows that any pair of distinct lines are incident at a single point.

Lemma 1.2 (Lines in the Projective Plane). *Given any two distinct lines in $K\mathbb{P}^2$, there is exactly one point incident to both of them.*

You will notice that this lemma is precisely the same as the previous, only with the words “point” and “line” interchanged. This is not a coincidence: points and lines in projective planes are *dual* in the sense that any result about points and lines in projective planes will still hold true if the two are interchanged.

This is also the reasoning for the choice of wording “incident to” for describing the relation between points and lines, rather than saying “two lines meet at a point” or “two points lie on a line”, since this makes the dualisation process easier.

1.2 Finite Projective Planes

Based on the previous algebraic construction, we define a combinatorial object.

A *finite projective plane (FPP)* is a finite set P of *points*, and a set $L \subseteq \mathcal{P}(P)$ of *lines* satisfying:

- (i) Every pair of points are incident to exactly one common line;
- (ii) Every pair of points are incident to exactly one common point;
- (iii) There are four points, no three of which belong to a single line.

The last condition is there only to rule out certain degenerate cases which lack the desired symmetries we like to work with.

Lemma 1.3 (Point-Line Matching). *Let (P, L) be an FPP, $\ell \in L$ be a line, and $p \in P$ be a point not incident to ℓ . Then, the number of points incident to ℓ is equal to the number of lines incident to p .*

Proof. By axiom (i), each point on ℓ is incident to exactly one line through p ; and by axiom (ii), each line through p is incident to exactly one point on ℓ . ■

Theorem 1.4 (FPP Structure). *Let (P, L) be an FPP. Then, there is a number q such that:*

- (i) Each line is incident to $q + 1$ points;
- (ii) Each point is incident to $q + 1$ lines;
- (iii) There are $q^2 + q + 1$ points;
- (iv) There are $q^2 + q + 1$ lines.

The number q is then called the *order* of the FPP.

Proof.

- (i) It suffices to show that any two lines are incident to the same number of points (and call this $q + 1$). Suppose ℓ and ℓ' are two lines. By Lemma 1.2, it is sufficient to find a point p not on either line, since each line would then be incident to as many points as there are lines incident to p .

Consider the 4 points p_1, p_2, p_3, p_4 guaranteed by axiom (iii). If one is in neither line, we are done. Otherwise, all four are on ℓ or ℓ' , and by axiom (iii), there must be two on each line, say, $p_1, p_2 \in \ell$ and $p_3, p_4 \in \ell'$. Now, consider the lines ℓ_{13} and ℓ_{24} connecting p_1 to p_3 and p_2 to p_4 , respectively. These lines meet at a point p .

If $p \in \ell$, then p_1 and p are points common to both ℓ and ℓ_{13} , so $\ell = \ell_{13}$ by uniqueness, and p_1, p_2, p_3 all lie on the line $\ell = \ell_{13}$, contradicting axiom (iii). Similarly, $p \notin \ell'$.

- (ii) Let p be a point. Again, by Lemma 1.1, it suffices to find a line not incident to p . Consider the 4 points p_1, p_2, p_3, p_4 guaranteed by axiom (iii), and without loss of generality, suppose $p \neq p_1$. Then, the line connecting p_1 and p_2 and the line connecting p_1 and p_3 cannot simultaneously contain p since they already both contain p_1 .
- (iii) Let p be a point and consider the $q + 1$ lines incident to it. Every pair of these lines intersect only at p , so each contains q points other than p , and the lines jointly cover the plane. So the total number of points is $(q + 1)q + 1 = q^2 + q + 1$.
- (iv) Each point is incident to $q + 1$ lines and every line is incident to $q + 1$ points, so the number of lines must be equal to the number of points.

■

1.3 Latin Squares

A *Latin square* is an $n \times n$ array of n distinct symbols such that every symbol appears in every row and every column.

Example.

A	B	A	B	C
		B	C	A
B	A	C	A	B

△

Example. Any Cayley table forms a Latin square. For instance, C_4 yields:

	0	1	2	3
0	0	1	2	3
1	1	2	3	0
2	2	3	0	1
3	3	0	1	2

△

Two $n \times n$ Latin squares $A = (a_{ij})$ and $B = (b_{ij})$ are *orthogonal* if the n^2 pairs (a_{ij}, b_{ij}) cover all possible pairs.

Example. If we pair

A	B	C		1	2	3
B	C	A	and	3	1	2
C	A	B		2	3	1

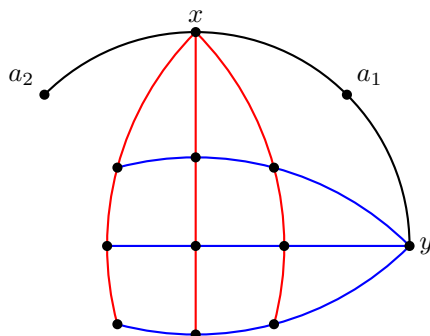
we obtain

(A,1) (B,2) (C,3)
 (B,3) (C,1) (A,2)
 (C,2) (A,3) (B,1)

All 9 possible pairs are present, so these Latin squares are orthogonal. △

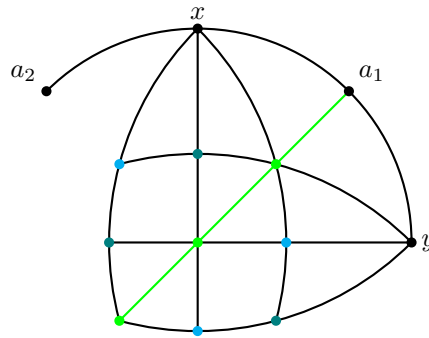
Let (P, L) be an FPP, and choose two points, say $x, y \in P$. They are connected by a line ℓ incident to $q + 1 - 2 = q - 1$ other points a_1, \dots, a_{q-1} , and $P \setminus \ell$ contains q^2 other points not on this line.

The point x is incident to q other lines, each disjointly incident to q points, so the x -lines hence partition these q^2 other points. Similarly, the point y is incident to q lines, each disjointly incident to q points, so the y -lines also partition these q^2 points. Also, each x -line meets all other y -lines, so the q^2 intersections have a Cartesian product structure and form a grid.



a_1 also lies on q other lines that also partition the grid. Also, each of these q lines meets each x -line and each y -line, so each line through a_1 is incident to q points within the grid; one in each row, and one in each column.

So, these q lines through a_1 generate a Latin square on the grid by labelling the points on the first line by a_{11} , points on the second line a_{12} , etc.



So in the example above, a_1 generates:

$$\begin{array}{ccc} a_{12} & a_{13} & a_{11} \\ a_{13} & a_{11} & a_{12} \\ a_{11} & a_{12} & a_{13} \end{array}$$

while the point a_2 generates:

$$\begin{array}{ccc} a_{21} & a_{23} & a_{22} \\ a_{22} & a_{21} & a_{23} \\ a_{23} & a_{22} & a_{21} \end{array}$$

The $q - 1$ points on the line connecting x and y thus generate $q - 1$ Latin squares.

Theorem 1.5. *The $q - 1$ Latin squares generated in this way are pairwise orthogonal.*

Proof. Without loss of generality, consider the Latin squares generated by the points a_1 and a_2 . These points each lie on q other lines corresponding to the symbols a_{1j} and a_{2j} used in their Latin squares, and each a_1 -line meets each a_2 -line at one of the q^2 points of the grid. So, every possible pair of symbols appears when the grids are merged. ■

Theorem 1.6. *There is an FPP of order $q > 1$ if and only if there are $q - 1$ pairwise orthogonal $q \times q$ Latin squares.*

2 Error-Correcting Codes

2.1 Introduction

Suppose Alice wishes to send Bob a message encoded in binary across an unreliable or *noisy* channel. That is, some bits in the message may be flipped during transmission.

One simple protocol to resist noise is for Alice to send each bit of the message repeatedly, say, ten times, then for Bob to take the most frequent received bit in each block of ten received bits to be the intended bit.

Message Bit	Code Bit
0	0000000000
1	1111111111

For instance, if Bob receives the string “10110111110100011000”, he can be fairly confident that the original message was “10”.

This replacement procedure constitutes an *error-correcting code*. The idea is that only certain strings of ten bits are valid or *admissible* strings, also called *codewords*, and that these admissible strings are selected to be very distinct from each other to minimise the chance that one is converted into by noise.

However, this code is not very *efficient*, because the rate of transmission decreases by a factor of ten when using this code.

Consider the similar repetition code:

Message Bit	Code Bit
0	00
1	11

This code detects one bit errors: if the string 01 is received, Bob will know there has been at least one bit flip. However, this code cannot detect two bit flips, as the intended message 00 could be converted into the admissible string 11 with two bit flips. The rate of this code is also 1/2.

If we use this code to send two bits of information, we have the encodings:

Message Bit	Code Bit
00	0000
01	0011
10	1100
11	1111

Again, this code is only safe against single bit flips. However, using three bits, we can achieve the same resilience against noise:

Message Bit	Code Bit
00	000
01	011
10	110
11	101

Any pair of these strings differ in two places, so again, two bits have to be corrupted to change one admissible string into another in this code. However, this code is also faster than the previous scheme, having rate 2/3.

Let the *alphabet* A be a set of symbols called *letters*. A *code* is a subset $C \subseteq A^n$, where n is the *length* of the code, and the elements of C are called *codewords*. If a code uses n bits of codewords to send k bits of plaintext, then the code has *rate* k/n .

Example. In the examples above, we have been working over the alphabet $A = \{0,1\}$, and we call such codes *binary*. △

An *encoding* is a bijection $e : W \rightarrow C$ from the set W of words in the plain text, to the code C .

Example. The table above describes an encoding from $W = \{0,1\}^2$ to a code $C = \{000,011,110,101\} \subseteq \{0,1\}^3$. △

For our purposes, it will not matter how this encoding is selected; all that is relevant is how “well-separated” the codewords are.

To quantify this separation, we define the *Hamming distance* of two codewords as the number of positions in which they differ. The Hamming distance forms a metric on any set of strings.

Example. The codewords 110 and 011 differ in the first and third positions, so they are Hamming distance $d(110, 011) = 2$ apart. \triangle

Note that the minimum separation $\min_{X \neq Y \in C} d(X, Y)$ of a code determines the maximum number of bit flips it can detect, as any number of bit flips exceeding this number could then potentially turn one codeword into another valid codeword. If for a code C , this minimum separation is

$$D := \min_{X \neq Y \in C} d(X, Y)$$

then we say that C is D -separated.

2.2 Block Codes

We will only be considering *block codes*, where the message is divided into blocks of fixed length k , each of which can be encoded without reference to any of the other blocks. That is, we will take the set of words W to be the set $\{0, 1\}^k$ of all possible binary strings of length k .

To *decode* a message encoded with a block code, we break the received transmission into blocks of length n , where n is the length of the code used. If a block is a codeword (i.e. admissible), then we assume that this block was correctly transmitted. Otherwise, we find the closest codeword to the received block, and interpret that as the intended codeword. If the code is well designed, this closest codeword should be unique.

Lemma 2.1. *If a block code is $(2r + 1)$ -separated, then it can correct r bit flips.*

Proof. An invalid block x can be at distance at most r from its nearest codeword y .

$$d(x, y) \leq r$$

Then, for any other codeword z , the reverse triangle inequality gives:

$$\begin{aligned} |d(z, y) - d(y, x)| &\leq d(x, z) \\ |(2r + 1) - r| &\leq d(x, z) \\ r + 1 &\leq d(x, z) \end{aligned}$$

so the codeword closest to x is unique. \blacksquare

Given a binary block code C , any encoding gives a bijection $W \rightarrow C$, we have $|C| = |W| = 2^k$, so the rate of a binary block code is simple to compute as:

$$\frac{\log_2 |C|}{n} = \frac{k}{n}$$

In a binary code, we can also interpret the alphabet $\{0, 1\}$ as the finite field $\mathbb{Z}_2 := \mathbb{Z}/2\mathbb{Z}$ of characteristic 2. If the codewords have length n , then they can be interpreted as elements of the vector space \mathbb{Z}_2^n .

In the code of length 3 above, we had the vectors

$$000, \quad 011, \quad 110, \quad 101$$

These elements form a linear subspace of \mathbb{Z}_2^3 , and we call this code a *linear code*.

The aim is to find codes in which every pair of codewords are far from each other in the Hamming metric. Linear codes have a simple feature that makes this easier to achieve:

Lemma 2.2. Suppose $C \subseteq \mathbb{Z}_2^n$ is a linear code such that every element of C other than $\mathbf{0}$ contains at least D coordinates equal to 1. Then, C is D -separated. That is,

$$\forall (X \neq Y \in C) : d(X, Y) \geq D$$

Proof. Suppose u and v are distinct codewords with $d(u, v) = r < d$. That is, they differ only in $r < d$ coordinates. Then, $u \oplus v$ is also a codeword, as C is linear. This codeword has a 1 only in the positions where u and v disagree, since vector addition in \mathbb{Z}_2^n is componentwise exclusive disjunction, so $u \oplus v$ has $r < d$ coordinates equal to 1, contradicting the construction of C . ■

It is also very simple to determine the rate of a linear code; if C is a linear subspace of \mathbb{Z}_2^n of dimension k , then it has 2^k elements, so the rate is just

$$\frac{\dim C}{n} = \frac{k}{n}$$

Example. Consider the code

$$C = \{000, 011, 110, 101\} \subset \mathbb{Z}_2^3$$

The minimum number of 1s in a non-zero codeword is 2, so C is 2-separated. C is also a 2-dimensional subspace of \mathbb{Z}_2^3 , so the rate is $\frac{2}{3}$. △

So, the goal is to look for k -dimensional subspaces of \mathbb{Z}_2^n whose non-zero elements contain a large number of 1s to get good separation. We also want k to be large to get a high rate. Since we want this subspace to have large dimension, it is usually easier to define it using $n - k$ linear equations, rather than using a basis of size k .

Example. The code above consists of the elements $(x_1, x_2, x_3) \in \mathbb{Z}_2^3$ which contain an even number of 1s, so they can be specified to be the elements satisfying

$$x_1 + x_2 + x_3 = 0$$

In other words, this code is defined by a *parity check*. △

2.3 Hamming Codes

In this section, we describe an efficient 3-separated binary code based on parity checks.

Let r be a positive integer and $n = 2^r - 1$ be the length of the code. The code will be an $(n - r)$ -dimensional subspace of \mathbb{Z}_2^n , so we need r linear equations to specify the codewords, and the rate will be almost 1:

$$\text{rate}(C) = \frac{\log_2(n)}{n} = \frac{n - r}{n} = 1 - \frac{r}{n} = 1 - \frac{r}{2^r - 1} \approx 1$$

for large r .

We arrange the linear equations into an $r \times n$ matrix B , so the code will be given by

$$C = \{x \in \mathbb{Z}_2^n : Bx = 0\}$$

The *Hamming code* of length n is given by the matrix with columns consisting of all binary numbers from 1 to n .

Example. For $r = 3$ and $n = 2^r - 1 = 7$, then B is given by

$$B = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{bmatrix}$$

△

Note that $(1,1,1,0,\dots,0) \in C$ regardless of r , since the first 3 columns of B have 0 below the top two rows. So, Hamming codes are at most 3-separated.

Lemma 2.3. *Hamming codes are precisely 3-separated.*

Proof. We have already shown that Hamming codes are at most 3-separated. If two words are less than distance 3 apart, then they must either be distance 1 or distance 2 apart. In the former case, both would fail to satisfy a parity check involving the different bit, and in the latter case, then there would be a non-codeword in between them, which would decode as either one of them. ■

As with any block code, to decode a message encoded with a Hamming code, we break the received code bits into blocks of length n , and check if the block is a codeword. If it is, then we assume that this block was correctly transmitted; otherwise, we find the closest codeword to the received block, and interpret that as the intended codeword.

We show that this closest codeword is at distance at most 1 for any possible received block, and furthermore, that this closest codeword is unique.

Theorem 2.4. *Let C be a Hamming code of length $n = 2^r - 1$. Then, for any invalid block $x \in \mathbb{Z}_2^n \setminus C$, there is a unique element y of C at distance $d(x,y) = 1$ from x .*

Proof. Suppose $x \in \mathbb{Z}_2^n \setminus C$ is not a codeword, so $Bx \neq \mathbf{0}$. Consider the vector

$$u := Bx \in \mathbb{Z}_2^r \setminus \{\mathbf{0}\}$$

By construction, the columns of B contain all possible non-zero vectors, so u must coincide with some column of B , say, the m th column B_m representing the binary number m .

Now, consider the vector

$$\tilde{x} := x \oplus \mathbf{e}_m$$

where \mathbf{e}_m is the standard m th basis vector, so \tilde{x} differs from x only in the m th coordinate. In particular, $d(x, \tilde{x}) = 1$.

Multiplying \tilde{x} by B , we have

$$\begin{aligned} B\tilde{x} &= Bx \oplus B\mathbf{e}_m \\ &= u \oplus B_m \\ &= B_m \oplus B_m \\ &= \mathbf{0} \end{aligned}$$

and we see that \tilde{x} is a codeword at distance 1 from x , as required.

For uniqueness, suppose u and v are distinct codewords at distance 1 from x . Then, by the triangle inequality, we have $d(u,v) \leq d(u,x) + d(x,v) = 2$, contradicting that Hamming codes are 3-separated. ■

Example. Decode the received string 1010011.

The first bit parity checks the positions whose unit digit is 1. That is, the 1st, 3rd, 5th, and 7th bits. We have $1 + 1 + 0 + 1 = 1$, so an error has occurred somewhere within these odd digits.

The second bit parity checks the positions whose 2s digit is 1. That is, the 2nd, 3rd, 6th, and 7th bits. We have $0 + 1 + 1 + 1 = 1$, so an error has occurred somewhere within these digits.

The fourth bit parity checks the positions whose 4s digit is 1. That is, the 4th, 5th, 6th, and 7th bits. We have $0 + 0 + 1 + 1 = 0$, so no error has occurred within these digits.

From this, we deduce that the 3rd digit has been flipped, so the closest codeword we correct to is 1000011. △

This process is somewhat involved, so we present a graphical method to quickly decide which positions to parity check:

Example. Decode the received string 1010011.

Arrange the string into a grid as follows, skipping the first entry:

	1	0	1
0	0	1	1

Now, perform parity checks within each of the highlighted regions:

	1	0	1
0	0	1	1

×

	1	0	1
0	0	1	1

×

	1	0	1
0	0	1	1

✓

The error is in the first two regions, so it must be in the 4th column. The error is not in the last region, so it must be in the first row.

	1	0	1
0	0	1	1

We deduce that the error is in the 3rd bit, so the closest codeword we correct to is 1000011. △

Example. Decode the received string 110101110100101.

Arrange the string into a grid as follows, skipping the first entry:

	1	1	0
1	0	1	1
1	0	1	0
0	1	0	1

Now, perform parity checks within each of the highlighted regions:

	1	1	0
1	0	1	1
1	0	1	0
0	1	0	1

✓

	1	1	0
1	0	1	1
1	0	1	0
0	1	0	1

×

	1	1	0
1	0	1	1
1	0	1	0
0	1	0	1

×

	1	1	0
1	0	1	1
1	0	1	0
0	1	0	1

✓

The error is in the second region, but not the first, so the error must be in the third column; the error is in the third region, but not the fourth, so the error is in the second row:

	1	1	0
1	0	1	1
1	0	1	0
0	1	0	1

So, we correct the received block to 110100110100101. △

Theorem 2.5 (Sphere-packing Bound). *Let C be a $(2r + 1)$ -separated binary code of length n . Then,*

$$|C| \sum_{i=0}^r \binom{n}{i} \leq 2^n$$

Proof. In Lemma 2.1, we showed that if a block code is $(2r + 1)$ -separated, then it can correct r bit flips, since every string that is distance at most r from a codeword x is closer to x than any other codeword. In other words, the balls of radius r centred on each codeword are all disjoint.

How many strings are contained in each ball of radius r ?

We count these strings based on their distance from the centre x . If we change 0 bits from x , then we just get the string x ; if we change 1 bit, then there are $n = \binom{n}{1}$ many strings at distance 1 from x ; if we change 2 bits, then there are $\binom{n}{2}$ many strings at distance 2 from x ; and so on, so the balls each contain

$$1 + \binom{n}{1} + \binom{n}{2} + \binom{n}{3} + \cdots + \binom{n}{r} = \sum_{i=0}^r \binom{n}{i}$$

strings. Since these balls are all disjoint, the total number of elements contained in all of these balls is $|C|$ times this sum, and there are 2^n possible strings, so

$$|C| \sum_{i=0}^r \binom{n}{i} \leq 2^n$$

as required. ■

A code that attains this bound is called a *perfect code*.

Lemma 2.6 (Domain of Codewords). *Let C be the Hamming code of length $n = 2^r - 1$. Then, each codeword is the closest codeword to n other elements of \mathbb{Z}_2^n .*

Proof. There are n possible bits to flip in each codeword. ■

Theorem 2.7. *Hamming codes are perfect codes.*

Proof. By the previous lemma, each codeword is closest to 2^r many possible bit strings: itself, and $n = 2^r - 1$ others adjacent to it. Also, there are $|C| = 2^{n-r}$ many codewords, so

$$\begin{aligned} |C| \cdot |B_r| &= 2^{n-r} 2^r \\ &= 2^n \end{aligned}$$
■

2.4 Shannon's Theorem

Suppose we have a binary communication channel which flips bits with probability p . We define the *Shannon capacity* to be

$$R := 1 + p \log_2(p) + (1 - p) \log_2(1 - p)$$

Theorem 2.8 (Shannon's Limit). *Using a binary communication channel which flips bits with probability p , there exists a code C with rate almost R and almost perfect accuracy.*

That is, for all $\varepsilon > 0$, there exists a code C with rate

$$\text{rate}(C) \geq 1 + p \log_2(p) + (1 - p) \log_2(1 - p) - \varepsilon$$

and such that the probability of decoding a codeword incorrectly is less than ε . Moreover, subject to any accuracy constraint, rates greater than R are not achievable.

Proof sketch. Choose a large value of n for the length of a block code and let F be a random variable measuring the number of bits flipped out of a message of length n . F has expectation $\mathbb{E}(F) = np$ and standard deviation $\sigma = \sqrt{npq}$ (where $q = 1 - p$). Notably, for large n , $\sigma \ll \mathbb{E}(F)$, so we will almost never have more than np bits flipped. Let $d := n(p + \varepsilon)$.

The first idea one might have is to find a $(2d + 1)$ -separated code with the given rate, but it turns out that this is extremely difficult to do.

Instead, choose M codewords from \mathbb{Z}_2^n uniformly and independently to form a code C . Now, suppose a codeword S is sent using this scheme, and is received as the string S' . There are two things that could go wrong during decoding:

A: More than d bits are flipped.

B: There is an incorrect codeword $Y \neq S$ at distance $d(S', Y) \leq d$ from S' .

The first case is rare by our choice of d , since $\sigma \ll \mathbb{E}(F) < d$. The second case occurs whenever one of the $M - 1$ codewords X in $C \setminus \{S\}$ is within distance d of S'

$$\mathbb{P}(B) = (M - 1) \frac{\# \text{ of strings } X \text{ with } d(S', X) \leq d}{2^n}$$

As before, the number of strings at distance r from S' is given by $\binom{n}{r}$, so the numerator is given by the sum

$$= (M - 1) \frac{\sum_{i=0}^d \binom{n}{i}}{2^n}$$

If d is not too large, then the sum of binomial coefficients is approximately the last summand, so $\mathbb{P}(B) \approx (M - 1) \frac{\binom{n}{d}}{2^n}$. So, we need $(M - 1) \binom{n}{d}$ to be small compared to 2^n :

$$(M - 1) \binom{n}{d} = \alpha 2^n$$

for some small constant $\alpha > 0$.

The rate of C is then given by

$$\begin{aligned} \text{rate}(C) &:= \frac{\log_2(M)}{n} \\ &\approx \frac{1}{n} \log_2 \left(\frac{2^n \alpha}{\binom{n}{d}} \right) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \left(\log_2(2^n) + \log_2(\alpha) - \log_2 \binom{n}{d} \right) \\
&= \frac{1}{n} \left(n + \log_2(\alpha) - \log_2 \binom{n}{d} \right) \\
&= 1 + \frac{1}{n} \log_2(\alpha) - \frac{1}{n} \log_2 \binom{n}{d}
\end{aligned}$$

$\frac{1}{n} \log_2(\alpha)$ is very small even for small α , so,

$$\begin{aligned}
&\approx 1 - \frac{1}{n} \log_2 \binom{n}{d} \\
&= 1 - \frac{1}{n} \log_2 \left(\frac{n!}{d!(n-d)!} \right) \\
&\approx 1 - \frac{1}{n} \log_2 \left(\frac{n!}{(np)!(n(1-p))!} \right)
\end{aligned}$$

By Stirling's formula and discarding sublinear factors,

$$\begin{aligned}
&\approx 1 - \frac{1}{n} \log_2 \left(\frac{n^n}{(np)^{np} (n(1-p))^{n(1-p)}} \right) \\
&= 1 - \frac{1}{n} \log_2 \left(\frac{n^n}{n^{np} p^{np} n^{n(1-p)} (1-p)^{n(1-p)}} \right) \\
&= 1 - \frac{1}{n} \log_2 \left(\frac{n^n}{n^n p^{np} (1-p)^{n(1-p)}} \right) \\
&= 1 - \frac{1}{n} \log_2 \left(\frac{1}{p^{np} (1-p)^{n(1-p)}} \right) \\
&= 1 - \frac{1}{n} \log_2 \left(\left(\frac{1}{p^p (1-p)^{(1-p)}} \right)^n \right) \\
&= 1 - \log_2 \left(\frac{1}{p^p (1-p)^{(1-p)}} \right) \\
&= 1 + \log_2 (p^p (1-p)^{(1-p)}) \\
&= 1 + p \log_2(p) + (1-p) \log_2(1-p)
\end{aligned}$$

■

The point is that, if we fix the rate, then the number d of bit flips we can correct is proportional to n , and if we fix the probability p , then the expected number of bits flipped is np , also proportional to n . However, the standard deviation $\sigma = \sqrt{npq}$ grows slower than proportionally to n , so as n increases, the chance that $n(p + \varepsilon)$ bits are flipped decreases to 0.

3 Discrete Geometry

A set $C \subseteq \mathbb{R}^d$ is *convex* if for all pairs $x, y \in C$, we have $\lambda x + (1 - \lambda)y \in C$ for all $\lambda \in [0, 1]$. That is, the line segment connecting x to y is also contained in C .

Example. The unit ball of any normed space is convex. △

Lemma 3.1. *The arbitrary intersection of convex sets is convex.*

Proof. Let $\{S_i\}_{i=1}^n$ be a family of convex sets. Then, for any $x, y \in \bigcap_{i=1}^n S_i$ we have $x, y \in S_i$ for all i , and all the S_i are convex, so $\lambda x + (1 - \lambda)y \in S_i$ for all i , so $\lambda x + (1 - \lambda)y \in \bigcap_{i=1}^n S_i$ and $\bigcap_{i=1}^n S_i$ is convex. ■

The expression $\lambda x + (1 - \lambda)y$ is called a *convex combination* of x and y . More generally, the convex combination of a collection of points $x_1, \dots, x_m \in \mathbb{R}^d$ is a point of the form

$$\sum_{i=1}^m \lambda_i x_i$$

where $\lambda_i \geq 0$, and $\sum_{i=1}^m \lambda_i = 1$.

We write $\text{cc}(E)$ to denote the set of all convex combinations of a set E .

Lemma 3.2. *The convex combinations operator is idempotent:*

$$\text{cc}(\text{cc}(E)) = \text{cc}(E)$$

Lemma 3.3. *A set E is convex if and only if it contains all of its convex combinations:*

$$E = \text{cc}(E)$$

Proof. For the forward direction, we induct on m . For $m = 1$, a convex combination of points in E is just a point in E . Now suppose E contains all convex combinations of at most m of its points. Then, we can reduce

$$\begin{aligned} \sum_{i=1}^{m+1} \lambda_i x_i &= \lambda_{m+1} x_{m+1} + \sum_{i=1}^m \lambda_i x_i \\ &= \lambda_{m+1} x_{m+1} + \left(\sum_{i=1}^m \lambda_i \right) \left(\sum_{i=1}^m \frac{\lambda_i}{\sum_{j=1}^m \lambda_j} x_i \right) \\ &= \lambda_{m+1} x_{m+1} + (1 - \lambda_{m+1}) \left(\sum_{i=1}^m \frac{\lambda_i}{\sum_{j=1}^m \lambda_j} x_i \right) \end{aligned}$$

The sum on the right is a convex combination, and by the inductive hypothesis, this is an element of E . Then, by convexity of E , the whole expression is a point in E .

For the reverse direction, if E contains all convex combinations of its points, then it contains all convex combinations of two of its points, which is the definition of convexity. ■

Given a set $E \subseteq \mathbb{R}^d$, we define the *convex hull* of E to be the intersection of all convex sets containing E .

$$\text{conv}(E) := \bigcap_{\substack{C \supseteq E \\ C \text{ convex}}} C$$

As the intersection of convex sets, the convex hull is convex. The convex hull also satisfies:

- (i) $E \subseteq \text{conv}(E)$, since every set in the intersection contains E ;
- (ii) If C is convex and $E \subseteq C$, then $\text{conv}(E) \subseteq C$, since the intersection is the minimal element of the poset of convex sets containing E .

Theorem 3.4. *For any $E \subseteq \mathbb{R}^d$,*

$$\text{conv}(E) = \text{cc}(E)$$

Proof. Note that $\text{cc}(E)$ is convex and that $E \subseteq \text{cc}(E)$ since every point in E is a convex combination of itself. So, by property (ii) of convex hulls, $\text{conv}(E) \subseteq \text{cc}(E)$.

Any point of $\text{cc}(E)$, i.e. a convex combination of points in E , is also a convex combination of points in $\text{conv}(E)$ since $E \subseteq \text{conv}(E)$. Because $\text{conv}(E)$ is convex, it contains all of these convex combinations, so $\text{cc}(E) \subseteq \text{conv}(E)$. ■

3.1 Separation

We will be concerned almost entirely with convex sets that are closed, and in almost all cases, they will also be bounded and hence compact (by Heine-Borel).

Lemma 3.5. *Every linear functional $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is of the form*

$$x \mapsto \langle x, y \rangle$$

where y is some fixed non-zero vector in \mathbb{R}^d .

Proof. Let ϕ be a linear functional. Define $y_i := \phi(\mathbf{e}_i)$ for each $0 \leq i \leq d$, where \mathbf{e}_i is the i th standard basis vector. Then, if $x = \sum_{i=1}^d x_i \mathbf{e}_i$, we have

$$\begin{aligned} \phi(x) &= \phi\left(\sum_{i=1}^d x_i \mathbf{e}_i\right) \\ &= \sum_{i=1}^d x_i \phi(\mathbf{e}_i) \\ &= \sum_{i=1}^d x_i y_i \\ &= \langle x, y \rangle \end{aligned}$$

■

We define a *hyperplane* in \mathbb{R}^d to be a set of the form

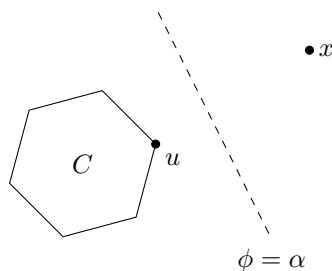
$$\Pi = \{x \in \mathbb{R}^d : \phi(x) = \alpha\}$$

for some non-zero linear functional ϕ and constant $\alpha \in \mathbb{R}$. Equivalently, it is an affine subspace of codimension 1.

Theorem 3.6 (Separation Principle I). *If $C \subseteq \mathbb{R}^d$ is compact and convex, and $x \in \mathbb{R}^d \setminus C$, then there is a hyperplane separating x from C . That is, there exists a linear functional $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ and a number α such that*

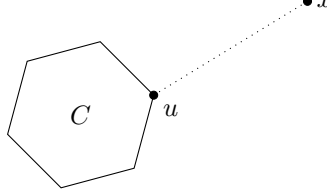
- $\phi(x) > \alpha$;
- $\phi(c) < \alpha$ for all $c \in C$.

Example.



△

Proof. Consider the function $C \rightarrow \mathbb{R}$ defined by $c \mapsto \|x - c\|$ that returns the distance from a point in C to the point x . This function is continuous, and so has a minimum on C . That is, there is a closest point u of C to x . Since $x \notin C$, $u \neq x$, so $\|x - c\| > 0$.



Now, define $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ by $y \mapsto \langle x - u, y \rangle$. We have

$$\begin{aligned} \phi(x) - \phi(u) &= \langle x - u, x \rangle - \langle x - u, u \rangle \\ &= \langle x - u, x - u \rangle \\ &= \|x - u\|^2 \\ &> 0 \end{aligned}$$

so $\phi(u) < \phi(x)$. Let α satisfy $\phi(u) < \alpha < \phi(x)$. To complete the proof, it remains to show that if $c \in C$, then $\phi(c) \leq \alpha$.

Suppose $c \in C$, but $\phi(c) > \alpha$. Consider the convex combination $p := \delta c + (1 - \delta)u$. Then,

$$\begin{aligned} \|x - p\|^2 &= \|x - \delta c - (1 - \delta)u\|^2 \\ &= \|(x - u) - \delta(c - u)\|^2 \\ &= \|x - u\|^2 - 2\delta\langle x - u, c - u \rangle + \delta^2\|c - u\|^2 \\ &= \|x - u\|^2 - 2\delta\phi(c - u) + \delta^2\|c - u\|^2 \end{aligned}$$

By assumption, $\phi(c - u) < 0$, and for small δ , $\delta^2\|c - u\|^2 \ll 2\delta\langle x - u, c - u \rangle$, so,

$$< \|x - u\|^2$$

so p is closer to x than u , contradicting the construction of u . ■

A *half-space* of \mathbb{R}^d is a set

$$H = \{x \in \mathbb{R}^d : \phi(x) \leq \alpha\}$$

for some non-zero linear functional ϕ and constant $\alpha \in \mathbb{R}$.

Corollary 3.6.1. *If $C \subset \mathbb{R}^d$ is a compact convex set, then C can be expressed as an intersection of half-spaces.*

Proof. For each point $x \in \mathbb{R}^d \setminus C$, there is a half-space containing C and not x , so the intersection of all half-spaces containing C will exclude all points $x \in \mathbb{R}^d \setminus C$ and hence is equal to C . ■

Given a set C , a *supporting hyperplane* of C is a hyperplane H that contains a boundary point x of C , but does not intersect the interior of C . Or equivalently, the (non-zero) linear functional ϕ given by the orthogonal vector of the hyperplane satisfies $\phi(c) \leq \phi(x)$ for all $c \in C$.

Theorem 3.7 (Supporting Hyperplanes). *If $C \subset \mathbb{R}^d$ is compact and convex, and $x \in \partial C$, then there is a hyperplane supporting C at x . That is, there is a non-zero linear functional $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\phi(c) \leq \phi(x)$ for all $c \in C$.*

Proof. Let $(x_i)_{i=1}^\infty \subseteq \mathbb{R}^d \setminus C$ be a sequence of points converging to x . By the separation principle, there exists, for each i , a linear functional ϕ_i defined by

$$u \mapsto \langle u, v_i \rangle$$

and a number α_i such that

- $\phi_i(x_i) > \alpha_i$;
- $\phi_i(c) < \alpha_i$ for all $c \in C$.

Without loss of generality, suppose each v_i is a unit vector, and hence that they have an accumulation point v which is also a unit vector. Passing to a subsequence and re-indexing, assume that $(v_i) \rightarrow v$.

For each i , we have

$$\langle x, v_i \rangle < \alpha_i < \langle x_i, v_i \rangle$$

Taking limits as $i \rightarrow \infty$, the inner products converge to $\langle x, v \rangle$, so $(\alpha_i) \rightarrow \langle x, v \rangle$. Then, for each $c \in C$, we have $\langle c, v_i \rangle \leq \alpha_i$, so taking limits, we have $\langle c, v \rangle \leq \langle x, v \rangle$, as required. ■

We can relax the hypotheses of the separation principle by only requiring that C is closed and not compact:

Theorem 3.8 (Separation Principle II). *If $C \subseteq \mathbb{R}^d$ is closed and convex, and $x \in \mathbb{R}^d \setminus C$, then there is a hyperplane separating x from C . That is, there exists a linear functional $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ and a number α such that*

- $\phi(x) > \alpha$;
- $\phi(c) < \alpha$ for each $c \in C$.

Proof. Let $c \in C$ and define $R := \|x - c\|$. Now, consider the intersection of C with the ball $B_R(x)$ of radius R centred on x . This is a compact set, so it has a point u closest to x . The proof from this point onwards is then identical to the first form of the theorem. ■

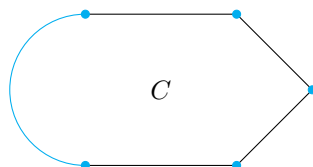
3.2 Extrema

Given a convex set $C \subseteq \mathbb{R}^d$, an *extreme point* of C is a point $c \in C$ not in the interior of any line segment contained in C . That is, if c is an extreme point and $x, y \in C$ satisfy

$$c = \lambda x + (1 - \lambda)y$$

with $\lambda \in (0, 1)$, then $x = y = c$.

Example. In the following figure, the blue points are extreme points.



△

Lemma 3.9 (Extreme Points of Faces). *Let H be a supporting hyperplane to the compact convex set C . Then,*

- (i) $H \cap C$ is compact and convex;
- (ii) Every extreme point of $H \cap C$ is an extreme point of C .

Proof.

- (i) The intersection of a compact and closed set is compact, and the intersection of convex sets is convex, so $H \cap C$ is compact and convex.
- (ii) Now, suppose x is an extreme point of $H \cap C$, but not an extreme point of C , so it is in the interior of a line segment in C . Because x is extreme in $H \cap C$, this line segment cannot be contained in $H \cap C$, so the segment must have endpoints on either side of H . But, this is impossible, since C is on one side of H .

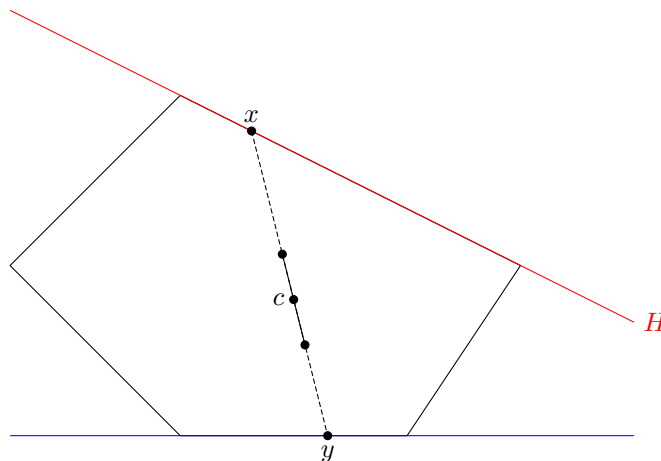
■

Theorem 3.10 (Extreme Point Theorem). *Let $C \subseteq \mathbb{R}^d$ be compact and convex, and let E be the set of its extreme points. Then,*

$$C = \text{conv}(E) = \text{cc}(E)$$

Proof. We induct on d . For $d = 1$, this is trivial.

We already know that $\text{conv}(E) \subseteq C$, so we show the other inclusion. Let $c \in C$. If c is an extreme point, then there is nothing to prove. So, supposing otherwise, c lies on a line segment in C , which we may extend in each direction until it intersects the boundary of C , at points, say, x and y .



Let H be a supporting hyperplane to C at x . Then, $H \cap C$ is a compact convex set by part (i) of the previous lemma, and it has codimension at least 1, so by the strong inductive hypothesis, $H \cap C$ is the convex hull of its extreme points, so

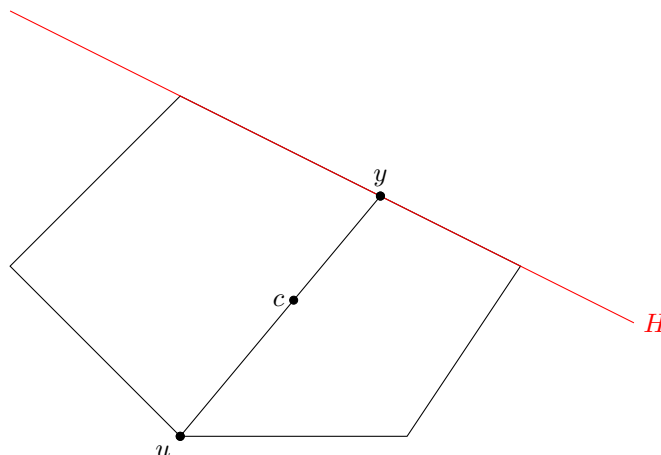
$$x \in H \cap C = \text{conv}(E_{H \cap C})$$

By part (ii) of the previous lemma, $E_{H \cap C} \subseteq E$, so $x \in \text{conv}(E)$. Through an identical argument, we have $y \in \text{conv}(E)$. Then, c lies on the line segment connecting x and y , so we also have $c \in \text{conv}(E)$, as required. ■

Theorem 3.11 (Caratheodory). *Each point of a compact convex set $C \subset \mathbb{R}^d$ is a convex combination of at most $d + 1$ of its extreme points.*

Proof. We induct on d . For $d = 1$, this is trivial.

Let $c \in C$ and choose an extreme point u of C . Consider the line passing through c and u . This line intersects the boundary of C at u on one side of c , and at a point y on the other.



Now, let H be a supporting hyperplane to C at y . By the strong inductive hypothesis, y is a convex combination of at most d extreme points in $H \cap C$, so c is a convex combination of these at most d points, and an extra point u . ■

3.3 Polyhedra and Polytopes

Here are two important constructions of convex sets in \mathbb{R}^d :

- A *polyhedron* is a bounded intersection of a finite set of half-spaces.
- A *polytope* is the convex hull of a finite set E .

Theorem 3.12 (Polyhedra are Polytopes). *Every polyhedron $C \subset \mathbb{R}^d$ is a polytope.*

Proof. Let

$$C := \bigcap_{i=1}^n S_i$$

be the bounded intersection of half-spaces S_i bounded by hyperplanes H_i . Let E be the set of extreme points of C .

We claim that every extreme point of C is the intersection of at most d hyperplanes. That is, for all $e \in E$, there exists $I \subseteq [n]$ with $|I| \leq d$ such that

$$\{e\} = \bigcap_{i \in I} H_i$$

This would imply that there are at most $\binom{n}{d}$ (i.e. finitely many) extreme points, so $C = \text{conv}(E)$ would be a polytope.

To prove the claim, we induct on d . For $d = 1$, this is trivial.

Let $e \in E$. If it were in the interior of all the half-spaces, then it would be in the interior of the intersection, C , in which case, e is not extreme. So, e must be on one of the hyperplanes, say H_1 .

Note that $H_1 \cap C = \bigcap_{i=1}^n (H_1 \cap S_i)$, and each $H_1 \cap S_i$ is a half-space, so $H_1 \cap C$ is a polyhedron of dimension at most $d - 1$. Moreover, e is extreme in $H_1 \cap C$, since if it were not, any line witnessing this would also witness this in C , contradicting that $e \in E$.

So, by the strong inductive hypothesis, in $H_1 \cap C$, we have

$$\{e\} = \bigcap_{i \in I'} (H_1 \cap H_i)$$

where $|I'| \leq d - 1$. So, in C , we have

$$\begin{aligned}\{e\} &= H_1 \cap \bigcap_{i \in I'} H_i \\ &= \bigcap_{i \in I' \cup \{1\}} H_i\end{aligned}$$

so e is the intersection of at most $|I' \cup \{1\}| \leq (d - 1) + 1 = d$ hyperplanes, completing the induction. This proves the claim, and the result follows. \blacksquare

3.4 Polars

Given a compact convex set $C \subseteq \mathbb{R}^d$, we define its *polar* to be the set

$$C^\circ := \{y \in \mathbb{R}^d : \forall x \in C, \langle x, y \rangle \leq 1\}$$

Under very weak conditions, polarity gives a bijection between C and C° :

Lemma 3.13. *If $C \subseteq \mathbb{R}^d$ is a compact convex set containing $\mathbf{0}$, then polarity is an involution:*

$$C^{\circ\circ} = C$$

Proof. By definition of a polar, for all $x \in C$ and $y \in C^\circ$, we have $\langle x, y \rangle = \langle y, x \rangle \leq 1$. By symmetry of the inner product, we have $\langle y, x \rangle \leq 1$, and

$$C^{\circ\circ} = \{x' \in \mathbb{R}^d : \forall y \in C^\circ, \langle y, x' \rangle \leq 1\}$$

so $C \subseteq C^{\circ\circ}$.

Now, suppose $C^{\circ\circ} \not\subseteq C$, so there exists $x \in C^{\circ\circ} \setminus C$ satisfying $\langle x, y \rangle \leq 1$ for all $y \in C^\circ$.

By the separation principle, there exists a linear functional $\phi(v) = \langle v, u \rangle$ for some fixed u and a constant α such that $\phi(x) > \alpha$ and $\phi(c) < \alpha$ for all $c \in C$.

Since $\mathbf{0} \in C$, $\alpha > \phi(0) = \langle 0, u \rangle > 0$, so by rescaling the orthogonal vector u to $u' := \frac{1}{\alpha}u$, we may assume that the constant is $\alpha' = 1$, so $\phi(c) = \langle c, u' \rangle < \alpha' = 1$ for all $c \in C$, and hence $u' \in C^\circ$. But then, $\phi(x) = \langle x, u' \rangle > \alpha' = 1$, so $x \notin C^{\circ\circ}$, contradicting our choice of x . \blacksquare

Lemma 3.14 (Polytope Polars). *If $C = \text{conv}(\{x_i\}_{i=1}^m)$, then*

$$C^\circ = \{y \in \mathbb{R}^d : \forall i, \langle x_i, y \rangle \leq 1\}$$

That is, we only have to check that $\langle x, y \rangle \leq 1$ for the vertices x_i , and not every point $x \in C$.

Proof. Define

$$C' := \{y \in \mathbb{R}^d : \forall i, \langle x_i, y \rangle \leq 1\}$$

Any $y \in C^\circ$ satisfies $\langle x_i, y \rangle \leq 1$ for all x_i , since $x_i \in C$, so $C^\circ \subseteq C'$.

For the reverse inclusion, let $y \in C'$. Then, any $x \in C = \text{conv}(\{x_i\}_{i=1}^m)$ is a convex combination

$$x = \sum_{i=1}^m \lambda_i x_i$$

so

$$\langle x, y \rangle = \sum_{i=1}^m \lambda_i \langle x_i, y \rangle$$

$$\begin{aligned} &\leq \sum_{i=1}^m \lambda_i \\ &= 1 \end{aligned}$$

so $y \in C^\circ$. Since y was arbitrary, we have $C' \subseteq C^\circ$. ■

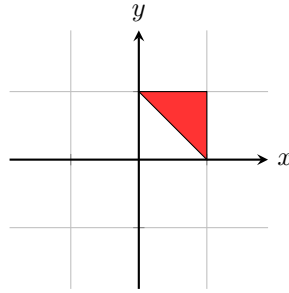
This lemma allows us to interpret the vectors of C as facets of C° . More precisely, notice that for any fixed $i \leq m$, the set

$$\{y \in \mathbb{R}^d : \langle x_i, y \rangle \leq 1\}$$

is a half-space (i.e. with orthogonal vector x_i and $\alpha = 1$), so this lemma equivalently says that C° is the intersection of these m half-spaces. From this, we deduce:

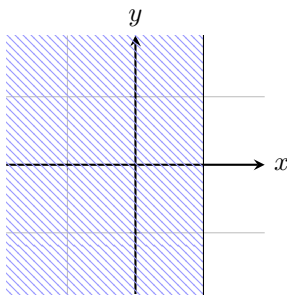
Corollary 3.14.1. *If C is a polytope, then C° is an intersection of half-spaces, and is hence a polyhedron if it is bounded.*

Example. Let $C = \text{conv}(\{(1,0), (0,1), (1,1)\})$:

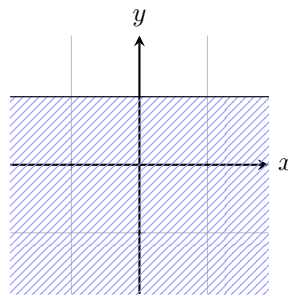


By the previous lemma, C° is the intersection of half-spaces

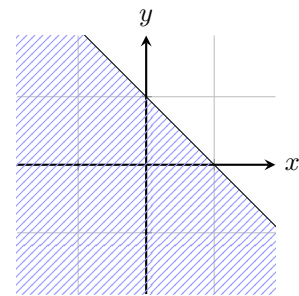
$$\begin{aligned} \{(x,y) : \langle (1,0), (x,y) \rangle \leq 1\} &= \{(x,y) : x \leq 1\} \\ \{(x,y) : \langle (0,1), (x,y) \rangle \leq 1\} &= \{(x,y) : y \leq 1\} \\ \{(x,y) : \langle (1,1), (x,y) \rangle \leq 1\} &= \{(x,y) : x + y \leq 1\} \end{aligned}$$



$$S_1 = \{(x,y) : x \leq 1\}$$

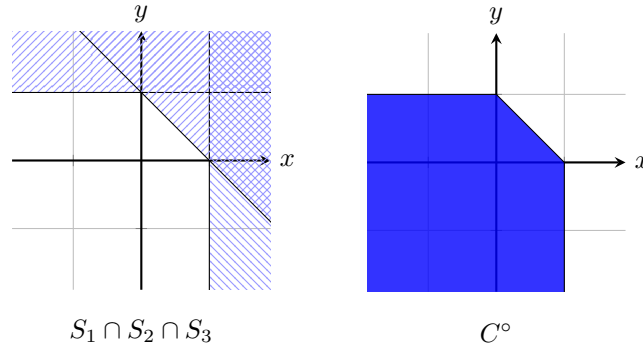


$$S_2 = \{(x,y) : y \leq 1\}$$



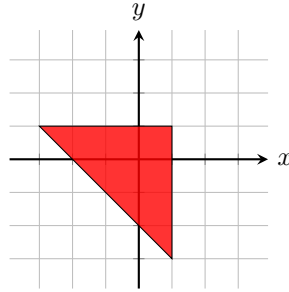
$$S_3 = \{(x,y) : x + y \leq 1\}$$

Shading the unwanted region so the intersection is easier to see, we have:



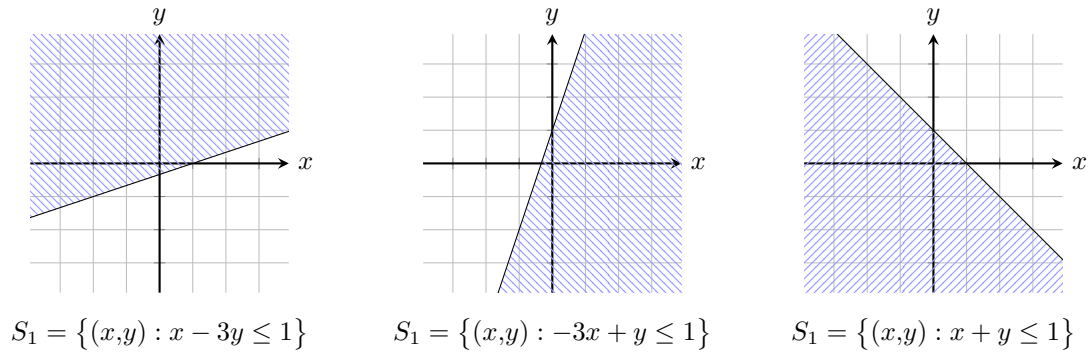
Note that in this case, C does not contain $\mathbf{0}$, and C° is unbounded, and hence not a polygon. \triangle

Example. Let $C = \text{conv}(\{(1, -3), (-3, 1), (1, 1)\})$:

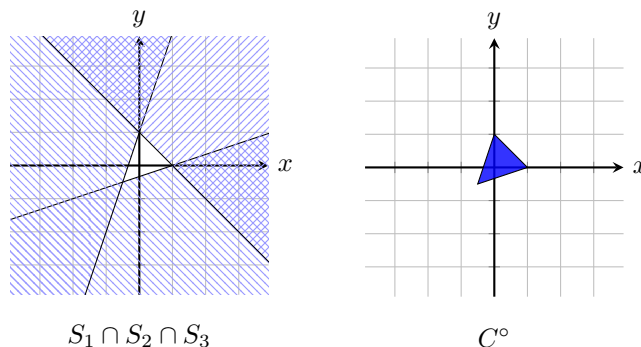


By the previous lemma, C° is the intersection of half-spaces

$$\begin{aligned} \{(x, y) : \langle (1, -3), (x, y) \rangle \leq 1\} &= \{(x, y) : x - 3y \leq 1\} \\ \{(x, y) : \langle (-3, 1), (x, y) \rangle \leq 1\} &= \{(x, y) : -3x + y \leq 1\} \\ \{(x, y) : \langle (1, 1), (x, y) \rangle \leq 1\} &= \{(x, y) : x + y \leq 1\} \end{aligned}$$



Shading the unwanted region so the intersection is easier to see, we have:



This time, $\mathbf{0} \in C$, so $C = C^{\circ\circ}$ and C° is a polytope whose vertices correspond to the facets of C . So, another way to find the vertices of the polar is to find the lines in which the facets lie.

In this case, we have the lines $1x + 0y = 1$, $0x + 1y = 1$, and $1x + 1y = -2$, which rearranges to $-\frac{1}{2}x - \frac{1}{2}y = 1$, so the vertices are $(1,0)$, $(0,1)$, and $(-\frac{1}{2}, -\frac{1}{2})$. \triangle

Lemma 3.15 (Inversion). *If C and D are convex sets with $D \subseteq C$, then $C^\circ \subseteq D^\circ$.*

Proof. If $y \in C^\circ$, then $\langle x, y \rangle \leq 1$ for all $x \in C$. We have $D \subseteq C$, so y also satisfies $\langle x, y \rangle \leq 1$ for all $x \in D \subseteq C$, so $y \in D^\circ$. \blacksquare

Theorem 3.16 (Polytopes are Polyhedra). *Every polytope $C \subset \mathbb{R}^d$ is a polyhedron.*

Proof. By translating if necessary, we may assume that C contains $\mathbf{0}$.

The strategy is to prove that the polar C° is bounded, and is hence a polyhedron. We have previously proved that polyhedra are polytopes, so C° is also a polytope. So, we may repeat this argument with C° replacing C , giving that $C^{\circ\circ}$ is also a polyhedron. Then, $C = C^{\circ\circ}$ is a polyhedron, as required.

We induct on d . If $d = 1$, this is trivial as polyhedra and polytopes are both just intervals.

Let C be the convex hull of finitely many points. If C is contained in a hyperplane, then the result immediately follows from the inductive hypothesis, so suppose otherwise.

Pick a point $u \in C$. Since C is d -dimensional, the set

$$\{v - u : v \in C\}$$

spans \mathbb{R}^d . Pick a basis $(v_i - u)_{i=1}^d$ consisting of vectors of this form.

The convex hull of the points u, v_1, \dots, v_d has non-empty interior since it contains a ball of radius $r > 0$ say around the barycentre

$$p := \frac{1}{d+1}(u + v_1 + \dots + v_d)$$

By translating the polytope C , suppose that $p = \mathbf{0}$ and that C is the convex hull of x_1, \dots, x_m .

We claim that C° is bounded.

Suppose that $y \in C^\circ$ has norm $k > 0$ and define the point $x := \frac{r}{k}y$. Then, x has norm

$$\|x\| = \left\| \frac{r}{k}y \right\| = \left| \frac{r}{k} \right| \|y\| = \frac{r}{k}k = r$$

so $x \in C$. So, by the definition of a polar, y must satisfy

$$\langle x, y \rangle \leq 1$$

$$\begin{aligned}
\langle \frac{r}{k}y, y \rangle &\leq 1 \\
\frac{r}{k} \langle y, y \rangle &\leq 1 \\
\frac{r}{k} \|y\|^2 &\leq 1 \\
\frac{r}{k} k^2 &\leq 1 \\
k &\leq \frac{1}{r}
\end{aligned}$$

So, $C^\circ \subseteq \mathbb{B}_{1/r}$ and is hence bounded. So, C° is a polyhedron.

Repeating this argument with C° replacing C , we have that $C^{\circ\circ} = C$ is a polyhedron, as required. ■

Along with Theorem 3.12, we have proved that polytopes and polyhedra in \mathbb{R}^d are equivalent.

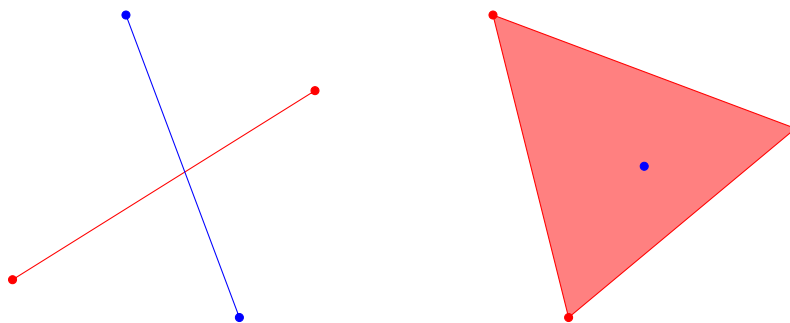
3.5 Radon's Lemma and Helly's Theorem

Lemma 3.17 (Radon). *Let $X \subseteq \mathbb{R}^d$ have cardinality $d+2$. Then, there exists a partition of X into two subsets whose convex hulls have non-empty intersection.*

Example. Consider $d = 2$, with 4 points in the plane. If at least 3 are colinear, then we can place the outermost points in one part, and the remaining points in the other part:



Otherwise, the points could form a convex quadrilateral, in which case, the diagonals intersect; alternatively, they could form a triangle, with a point inside:



△

Proof. Let $I = \{1, \dots, d+2\}$. Let the points be $(x_i)_{i \in I} \subseteq \mathbb{R}^d$. Adjoin an extra coordinate to each x_i , and set the coordinate equal to 1 to obtain $d+2$ points $(y_i)_{i \in I} \subseteq \mathbb{R}^{d+1}$.

Because we have $d+2$ points in \mathbb{R}^{d+1} , the y_i are not linearly independent, so there are scalars α_i not all zero such that

$$\sum_{i \in I} \alpha_i y_i = \mathbf{0}$$

Define the sets

$$A := \{i : \alpha_i > 0\}, \quad B := \{i : \alpha_i < 0\}$$

and define the scalars $\beta_i := -\alpha_i$, positive for $i \in B$. Sorting positive and negative coefficients, we have

$$\begin{aligned} \sum_{\substack{i \in I \\ \alpha_i > 0}} \alpha_i y_i + \sum_{\substack{i \in I \\ \alpha_i < 0}} \alpha_i y_i &= \mathbf{0} \\ \sum_{\substack{i \in I \\ \alpha_i > 0}} \alpha_i y_i &= \sum_{\substack{i \in I \\ \alpha_i < 0}} -\alpha_i y_i \\ \sum_{i \in A} \alpha_i y_i &= \sum_{i \in B} \beta_i y_i \end{aligned}$$

Considering only the first d coordinates, we have

$$\sum_{i \in A} \alpha_i x_i = \sum_{i \in B} \beta_i x_i$$

and considering only the final coordinate, we have

$$\sum_{i \in A} \alpha_i = \sum_{i \in B} \beta_i =: C$$

Because not all of the scalars are zero, A and B are non-empty, so these sums are non-empty, and we have $C > 0$.

Then, the vectors

$$x := \sum_{i \in A} \frac{\alpha_i}{C} x_i, \quad y := \sum_{i \in B} \frac{\beta_i}{C} x_i$$

are convex combinations of two disjoint subsets of the x_i , and these two vectors are equal. ■

Instead of partitioning X into two parts whose convex hulls have non-empty intersection, one could ask if we can partition X into three or more subsets whose convex hulls have non-empty intersection. For three parts, we clearly need to start with more than $d + 2$ points.

The right number turns out to be $2d + 3$, and in general, if we are partitioning into k subsets, we need $(k - 1)(d + 1) + 1$ points. This more general theorem holds, but is much harder to prove than Radon's lemma.

Theorem 3.18 (Tverberg). *Each set of $(k - 1)(d + 1) + 1$ points in \mathbb{R}^d can be partitioned into k subsets whose convex hulls have non-empty intersection.*

The simplest proof of Tverberg's theorem relies on the following result:

Theorem 3.19 (Colourful Caratheodory Theorem). *Let C_1, \dots, C_{d+1} be arbitrary subsets of \mathbb{R}^d , each coloured with a different colour. Suppose that $\mathbf{0} \in \text{conv}(C_i)$ for all $1 \leq i \leq d + 1$. Then, there is a rainbow set $R \subseteq \bigcup_i C_i$ with precisely one point of each colour whose convex hull contains $\mathbf{0}$.*

The next theorem we will prove is a striking dual of Caratheodory's theorem.

Recall that a space X is compact if and only if for every open cover \mathcal{U} of X , there exists a finite subcover $\mathcal{U}_0 \subseteq \mathcal{U}$:

$$\text{compact}(X) \quad \equiv \quad \bigcup \mathcal{U} = X \rightarrow \exists \text{ finite } \mathcal{U}_0 \subseteq \mathcal{U} : \bigcup \mathcal{U}_0 = X$$

We can De Morgan-dualise this statement by replacing open sets by closed sets, unions with intersections, and covers with empty-intersections:

$$\begin{aligned} &\equiv \quad \bigcap \mathcal{F} = \emptyset \rightarrow \exists \text{ finite } \mathcal{F}_0 \subseteq \mathcal{F} : \bigcap \mathcal{F}_0 = \emptyset \\ &\equiv \quad \bigcap \mathcal{F} \neq \emptyset \leftarrow \forall \text{ finite } \mathcal{F}_0 \subseteq \mathcal{F} : \bigcap \mathcal{F}_0 \neq \emptyset \end{aligned}$$

Taking the contrapositive in the second line, we have that a set is compact if and only if for every (non-empty) family \mathcal{F} of closed subsets of X , every finite subfamily $\mathcal{F}_0 \subseteq \mathcal{F}$ having non-empty intersection implies that \mathcal{F} has non-empty intersection.

The next theorem shows that in the convex case, we do not need to check all finite subfamilies, but only those of cardinality at most $d + 1$.

Theorem 3.20 (Helly's Theorem). *Let $\mathcal{F} = (C_i)_{i=1}^m$ be a family of convex sets in \mathbb{R}^d , and suppose that every subfamily $\mathcal{F}_0 \subseteq \mathcal{F}$ of cardinality at most $d + 1$ has non-empty intersection. Then, the whole family has a non-empty intersection:*

$$\bigcap \mathcal{F} \neq \emptyset$$

Proof. We induct on the number of sets m .

First, note that if $m \leq d + 1$, there is nothing to prove, since the desired result is included in the hypotheses of the theorem.

For the base case, suppose $m = d + 2$.

Then, for each C_i , the other $d + 1$ sets $C_{k \neq i}$ have non-empty intersection by assumption, so we may select a point x_i in each of these intersections:

$$x_i \in \bigcap_{k \neq i} C_k$$

By Radon's lemma, $X := \{x_i\}_{i=1}^{d+2}$ has a partition into two subsets $X_1, X_2 \subseteq X$ whose convex hulls have non-empty intersection. Let u be a point in this intersection.

$$u \in \text{conv}(X_1) \cap \text{conv}(X_2)$$

We claim that u is contained in each C_i , and is hence in $\bigcap \mathcal{F}$.

Fix some $1 \leq j \leq m$, and without loss of generality, suppose that x_j is in the Radon subset X_1 , so $x_j \notin X_2$. By construction of the x_i , we have $x_i \in C_j$ for all $i \neq j$, so $X_2 \subseteq X \setminus \{x_j\} \subseteq C_j$. Since C_j is convex, it also contains $\text{conv}(X_2) \ni u$, so $u \in C_j$.

In the above, we have assumed that the x_i are all distinct. But if this were not the case, say $x_i = x_j$ for some $i \neq j$, then by construction, $x_i \in C_k$ for all $k \neq i$, but also $x_i = x_j \in C_i$ (since $i \neq j$), so $x_i \in C_k$ for all k , and the intersection $\bigcap \mathcal{F}$ is again non-empty.

For the inductive step, suppose $m > d + 2$ and that the result holds for $m - 1$. Consider a new family of $m - 1$ sets given by

$$\mathcal{F}' := \{C_1 \cap C_2, C_3, C_4, \dots, C_m\}$$

and let $\mathcal{F}_0 \subseteq \mathcal{F}'$ be a subfamily of cardinality $d + 1$. Then, $\bigcap \mathcal{F}_0$ is the intersection of at most $d + 2$ of the original C_i , which is non-empty by the base case.

So, \mathcal{F} satisfies the hypotheses of the result, so by the induction hypothesis, $\bigcap \mathcal{F}' \neq \emptyset$. Then,

$$\emptyset \neq \bigcap \mathcal{F}' = (C_1 \cap C_2) \cap C_3 \cap \dots \cap C_m = \bigcap \mathcal{F}$$

which completes the inductive step. ■

If we also require that the C_i are compact, then Helly's theorem also holds for arbitrary collections \mathcal{F} , and not just finite collections.

4 Partially Ordered Sets and Set Systems

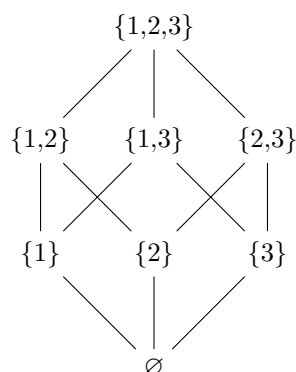
A relation \leq on a set X is a (*weak* or *non-strict*) *partial order* on X if it satisfies, for all $x, y, z \in X$:

- (i) reflexivity: $x \leq x$;
- (ii) transitivity: $(x \leq y \wedge y \leq z) \rightarrow x \leq z$;
- (iii) antisymmetry: $(x \leq y \wedge y \leq x) \rightarrow x = y$.

The pair (X, \leq) is then called a *partially ordered set* or *poset*.

Note that not all elements in a poset may be *comparable* under the ordering:

Example. Consider $(\mathcal{P}([3]), \subseteq)$, illustrated as the *Hasse diagram* below, where an edge from a vertex x travelling upwards to a vertex y indicates that $x \leq y$.



In this example, $\{1\}$ and $\{2,3\}$ are *incomparable* in this poset, because neither $\{1\} \subseteq \{2,3\}$ nor $\{2,3\} \subseteq \{1\}$ hold. \triangle

Vertices on the same horizontal level in a Hasse diagram are always incomparable.

If every pair of elements *are* comparable, then the ordering is *total*.

Example. (\mathbb{R}, \leq) is a total ordering. \triangle

Let \leq be a partial order on a set X .

A *chain* is a subset $C \subseteq S$ such that \leq is total on C . That is, every pair of elements in C are comparable under \leq :

$$\forall c_1, c_2 \in C : c_1 \leq c_2 \vee c_2 \leq c_1$$

Example. A chain in $(\mathcal{P}([3]), \subseteq)$ is given by the sequence of elements

$$\emptyset \subseteq \{1\} \subseteq \{1,2\} \subseteq \{1,2,3\}$$

or

$$\{3\} \subseteq \{1,3\}$$

\triangle

An *antichain* is a subset $A \subseteq S$ such that every pair of elements in A are incomparable under \leq :

$$\forall a_1, a_2 \in A : a_1 \not\leq a_2 \wedge a_2 \not\leq a_1$$

Example. An antichain in $(\mathcal{P}([3]), \subseteq)$ is given by the set

$$\{\{2\}, \{1,3\}\}$$

or

$$\{\{1\}, \{2\}, \{3\}\}$$

△

More generally, in $(\mathcal{P}([n]), \subseteq)$, any collection of subsets of a fixed cardinality form an antichain, since if two different sets have the same number of elements, then neither can be a subset of the other.

There are $\binom{n}{k}$ many subsets of $[n]$ of cardinality k , and this number is maximised when $k \approx n/2$. If n is even, then this has size precisely

$$\binom{n}{n/2}$$

and if n is odd, the two largest antichains are at level $k = (n-1)/2$ and $k = (n+1)/2$, and in either case, this has size

$$\binom{n}{\lfloor n/2 \rfloor}$$

Theorem 4.1 (Sperner). *The largest antichain in $(\mathcal{P}([n]), \subseteq)$ has cardinality*

$$\binom{n}{\lfloor n/2 \rfloor}$$

We will deduce Sperner's theorem from a stronger statement due to Lubell, Yamamoto, and Meshalkin:

Theorem 4.2 (LYM Inequality). *Let \mathcal{F} be an antichain in $(\mathcal{P}([n]), \subseteq)$. Then,*

$$\sum_{A \in \mathcal{F}} \frac{1}{\binom{n}{|A|}} \leq 1$$

Proof of Sperner's Theorem. It is clear that such an antichain exists – just pick all subsets with $\lfloor n/2 \rfloor$ elements.

To show that such an antichain is maximal, let \mathcal{F} be an antichain in $(\mathcal{P}([n]), \subseteq)$. Because $\binom{n}{\lfloor n/2 \rfloor} \geq \binom{n}{|A|}$ for any A , we have from the LYM inequality,

$$\begin{aligned} \sum_{A \in \mathcal{F}} \frac{1}{\binom{n}{|A|}} &\leq 1 \\ \sum_{A \in \mathcal{F}} \frac{1}{\binom{n}{\lfloor n/2 \rfloor}} &\leq 1 \\ |\mathcal{F}| \frac{1}{\binom{n}{\lfloor n/2 \rfloor}} &\leq 1 \\ |\mathcal{F}| &\leq \binom{n}{\lfloor n/2 \rfloor} \end{aligned}$$

■

Proof of the LYM Inequality. There is a bijection from the set of permutations on $[n]$ to the set of maximal chains in $(\mathcal{P}([n]), \subseteq)$, where each permutation gives the order in which to add elements to subsets in the chain.

For instance, for $n = 5$, the permutation $[5, 2, 3, 1, 4]$ corresponds to the chain

$$\emptyset \subseteq \{5\} \subseteq \{5, 2\} \subseteq \{5, 2, 3\} \subseteq \{5, 2, 3, 1\} \subseteq \{5, 2, 3, 1, 4\}$$

We pick a maximal chain/permutation R uniformly at random. The probability that R is any given maximal chain C is

$$\mathbb{P}(R = C) = \frac{1}{n!}$$

since there are $n!$ permutations on $[n]$.

For each subset $A \subseteq [n]$, let E_A be the set of maximal chains that contain A . If A and B are in both in \mathcal{F} , then E_A and E_B must be disjoint, since A and B must be incomparable and cannot both belong to the same chain. So,

$$\begin{aligned} \mathbb{P}(R \cap \mathcal{F} \neq \emptyset) &= \mathbb{P}\left(R \in \bigsqcup_{A \in \mathcal{F}} E_A\right) \\ &= \sum_{A \in \mathcal{F}} \mathbb{P}(R \in E_A) \end{aligned}$$

and since this is a probability, it is bounded above by 1:

$$\sum_{A \in \mathcal{F}} \mathbb{P}(R \in E_A) \leq 1$$

To finish the proof, it remains to show that

$$\mathbb{P}(R \in E_A) = \frac{1}{\binom{n}{|A|}}$$

First, we have

$$\mathbb{P}(R \in E_A) = \frac{\# \text{ of maximal chains containing } A}{\# \text{ of all maximal chains}}$$

A maximal chain containing A corresponds to a permutation which has the elements of A in any order as a prefix. For instance, any permutation starting with the numbers 1, 2, and 3 in any order will generate the subset $\{1,2,3\}$ in the corresponding chain.

So, there are $|A|!$ many ways to arrange this prefix, then $(n - |A|)!$ ways to arrange the remaining numbers. So,

$$\mathbb{P}(R \in E_A) = \frac{\# \text{ of maximal chains containing } A}{\# \text{ of all maximal chains}} = \frac{|A|!(n - |A|)!}{n!} = \frac{1}{\binom{n}{|A|}}$$

as required. ■

4.1 Dilworth's Theorem

A chain and an antichain can have at most one common element, for if x and y were two common elements, then the chain would require x and y to be comparable, and the antichain would require x and y to be incomparable.

So, if a poset can be covered with m chains, then there cannot be any antichains with more than m elements by the pigeonhole principle. Hence, another method of proving Sperner's theorem would be to find a cover of $(\mathcal{P}([n]), \subseteq)$ using

$$\binom{n}{\lfloor n/2 \rfloor}$$

chains. In fact, the proof of the LYM inequality above uses a similar, but simpler, idea, since we only looked at the covering generated by all maximal chains, and then counted how many times each set was covered.

This process also works in reverse to deduce that there is a covering by a small number of chains if there are only small antichains:

Theorem 4.3 (Dilworth). *Let (Ω, \leq) be a poset in which every antichain has at most m elements. Then, Ω can be covered by m chains (or fewer).*

Proof. We prove the case for finite Ω only.

We induct on $|\Omega|$. If $|\Omega| = 1$, then there is nothing to prove.

Suppose $|\Omega| > 1$, and that the result holds for all smaller posets. Let m be the size of the largest antichain in Ω , and choose a maximal chain $C = \{c_1 \leq c_2 \leq \dots \leq c_n\}$ in Ω .

Suppose $\Omega \setminus C$ has no antichains of length m , so every antichain has at most $m - 1$ elements. Then, the smaller poset $\Omega \setminus C$ may be covered by $m - 1$ chains (or fewer) by the inductive hypothesis, so Ω may be covered by m chains by adding C to this cover, and we are done.

Otherwise, $\Omega \setminus C$ has (maximal) antichains of length m . Let $A = \{a_1, \dots, a_m\} \subseteq \Omega \setminus C$ be such an antichain, and define the sets

$$A^- = \{x \in \Omega : \exists i, x \leq a_i\}, \quad A^+ = \{x \in \Omega : \exists i, x \geq a_i\}$$

Note that these sets jointly cover Ω , since if there were an $x \in \Omega$ but not $A^- \cup A^+$, then this x would be incomparable to all the a_i , so we could extend A , contradicting the maximality of A .

However, we also have that neither of these sets can be all of Ω , since if $C \subseteq A^-$, then we would have $c_n \leq a_i$ for some i , so we could extend C by adding a_i , contradicting the maximality of C ; and similarly, if $C \subseteq A^+$, we would have $c_1 \geq a_i$ for some i , again contradicting the maximality of C .

So, because A^- and A^+ are strict subsets of Ω , they are smaller posets, so the inductive hypothesis applies. So, A^- and A^+ can each be covered with m chains C_i^- and C_i^+ , respectively:

$$A^- = \bigcup_{i=1}^m C_i^- \quad A^+ = \bigcup_{i=1}^m C_i^+$$

Each of these decompositions partition the a_i , so, reindexing if necessary, we may assume that $a_i \in C_i^-$ and $a_i \in C_i^+$ for each i .

We claim that a_i is the maximal element of C_i^- and the minimal element of C_i^+ .

Suppose otherwise, so there is an element x such that $a_i < x \in C_i^-$. Since $x \in C_i^- \subseteq A^-$, we have $x \leq a_j \neq a_i$. But then, $a_i < a_j$ are comparable, contradicting that A is an antichain. The proof that a_i is minimal in C_i^+ is entirely symmetric.

Then, the m chains $C_i^- \cup C_i^+$ cover Ω . ■

4.2 Covering by Chains

Along with our previous observation that there cannot be an antichain longer than the number of chains in a covering, Dilworth's theorem gives that:

Corollary 4.3.1. *For any poset (Ω, \leq) , the length of the largest antichain is equal to the minimum number of chains required to cover Ω .*

We can also apply Dilworth's theorem to $(\mathcal{P}([n]), \subseteq)$ with $m = \binom{n}{\lfloor n/2 \rfloor}$ given by Sperner's theorem to obtain:

Corollary 4.3.2. *The poset $(\mathcal{P}([n]), \subseteq)$ can be covered using $\binom{n}{\lfloor n/2 \rfloor}$ chains.*

We can also show this directly, and conversely use this result with Dilworth's theorem to give another proof of Sperner's theorem. This proof will depend on Hall's theorem:

Theorem 4.4 (Hall). *Let $G = (L \cup R, E)$ be a bipartite graph. For each subset $U \subseteq L$, let $N_G(U)$ denote the open neighbourhood of U in G :*

$$N_G(U) := \{v \in R : \exists u \in U, (u, v) \in E\}$$

That is, the set of vertices in R that are adjacent to at least one element in U .

Then, there is an matching that covers L if and only if for every $U \subseteq L$,

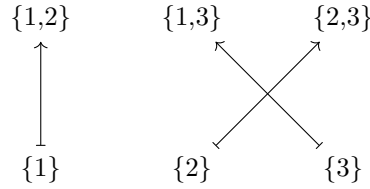
$$|U| \leq |N_G(U)|$$

That is, every subset $U \subseteq L$ must have sufficiently many neighbours in R for such a matching to exist.

Proof of Corollary 4.3.2. Let $r < n/2$, and consider the set R of subsets of cardinality r , and the set R^+ of subsets of cardinality $r + 1$:

$$R_r := \{S \subseteq [n] : |S| = r\}, \quad R_r^+ := \{S \subseteq [n] : |S| = r + 1\}$$

We claim that there is an injection $f : R_r \rightarrow R_r^+$ such that $A \subseteq f(A)$ for all $A \in R_r$. For instance,



Consider the bipartite graph $G_r = (R_r \cup R_r^+, E)$, where $(A, B) \in E$ if and only if $A \subseteq B$ (i.e. the subgraph of the Hasse diagram induced by taking the r and $(r + 1)$ th rows). Note that an R_r -saturated matching in this bipartite graph precisely corresponds to the required injection.

Each set $S \in R_r$ has $n - r$ neighbours, since r of the n numbers we could add are already in the set. Conversely, each set $S^+ \in R_r^+$ has $r + 1$ neighbours, since we may remove any of its $r + 1$ elements.

Now, let $U \subseteq R_r$. Then, there are

$$\sum_{S \in U} \deg(S) = |U|(n - r)$$

edges incident to U , and each vertex in R_r^+ is incident to at most $r + 1$ of these edges, so

$$|N_G(U)| \geq \frac{|U|(n - r)}{r + 1}$$

Since $r < n/2$, we have

$$\begin{aligned} r + 1 &\leq \frac{n}{2} \\ r + 1 &\leq n - \frac{n}{2} \\ r + 1 &\leq n - r \end{aligned}$$

so

$$|N_G(U)| \geq |U|$$

so Hall's condition is satisfied, and there exists an R_r -saturated matching M_r .

Repeating this construction on every layer, we obtain matchings from each layer to the next, up to layer $k := \lfloor n/2 \rfloor$.

Note that $\bigcup_{r=0}^k M_r$ is a subgraph of $\bigcup_{r=0}^k G_r$ consisting of disjoint paths: there can be no vertices of degree 3 or higher since every layer consists of perfect matchings. By the construction of the edge set of the G_r , each such path defines a chain.

We can cover the rest of $\mathcal{P}([n])$ by mirroring this construction as follows. Given $A \subseteq [n]$ with $|A| < k$, let $g(A) = [n] \setminus A$. Then, each above chain defines a chain $g(C)$.

If n is odd, then the middle two layers R_k and R_{k+1} have the same cardinality, and Hall's theorem gives a perfect matching, so the chains join up as desired. If n is even, then we may have some sets at $k = n/2$ uncovered, in which case, we may cover them with additional singleton chains.

Each chain contains a set in the middle layer k , so there are $\binom{n}{k} = \binom{n}{\lfloor n/2 \rfloor}$ total chains in this cover. ■

With the earlier observation that that a covering with m chains implies that every antichain has at most m elements, this provides another proof of Sperner's theorem.

There is a nice variation on this construction in which each chain is "symmetric": each chain consists of sets of sizes $0, 1, \dots, n-1, n$ or of sizes $1, 2, \dots, n-2, n-1$, etc. This improvement cannot be deduced from just Hall's theorem.

Theorem 4.5 (de Bruijn, Tengbergen, Kruyswijk). *The poset $(\mathcal{P}([n]), \subseteq)$ can be covered using $\binom{n}{\lfloor n/2 \rfloor}$ symmetric chains.*

Proof. We induct on n . For $n = 1$, there is the unique chain $(\emptyset, \{1\})$, which is symmetric.

Suppose $n > 1$, and that the result holds for $n-1$, so there is a decomposition of $[n-1]$ into symmetric chains. For each chain $C_i = (A_1, A_2, \dots, A_k)$ in this decomposition, we form two new chains

$$C_i^+ := (A_1, A_2, \dots, A_k, A_k \cup \{n\})$$

and

$$C_i^- := (A_1 \cup \{n\}, A_2 \cup \{n\}, \dots, A_{k-1} \cup \{n\})$$

The collection of these new chains covers $\mathcal{P}([n])$ since for each old subset $A \in \mathcal{P}([n-1])$, the new chains C_i^+ and C_i^- cover A and $A \cup \{n\}$, respectively.

If C_i consists of sets of size j to $n-j$, then C_i^+ has sizes j to $n-j+1 = (n+1)-j$, and C_i^- has sizes $j+1$ to $n-j = (n+1)-(j+1)$, so the new chains are also symmetric.

This construction also ensures that in each chain, each set has precisely one more element than the set before it, since the first transformation adds a new set one element larger to the top of an existing chain, and the second transformation adds one element to every set in an existing chain. In either case, this relation is preserved. So, each chain contains an element of cardinality $\lfloor n/2 \rfloor$, so there are $\binom{n}{\lfloor n/2 \rfloor}$ chains in the cover. ■

4.3 VC Dimension and the Sauer-Shelah Lemma

Given a finite set $U = \{u_1, u_2, \dots, u_m\}$, a family of sets \mathcal{F} *shatters* U if for every subset $V \subseteq U$, there is an element $A \in \mathcal{F}$ such that $A \cap U = V$.

Example. If U consists of the three vertices of a triangle in \mathbb{R}^2 , and \mathcal{F} is the family of half-spaces in \mathbb{R}^2 , then each of the 8 subsets of U can be obtained by intersecting U with an appropriate half-space, so \mathcal{F} shatters U .

However, if U consists of four points in the plane, then \mathcal{F} cannot shatter U by Radon's lemma. △

Given a set Ω and a family of sets $\mathcal{F} \subseteq \mathcal{P}(\Omega)$, we define the *Vapnik-Cervonenkis (VC) dimension* $\text{VC}(\mathcal{F})$ of \mathcal{F} to be maximum cardinality of a subset of Ω that \mathcal{F} can shatter.

Example.

- $\text{VC}(\mathcal{P}([n])) = n$;
- $\text{VC}(\{\text{half-spaces in } \mathbb{R}^n\}) = n + 1$;
- $\text{VC}(\text{any FPP}) = 2$.

△

Fix integers n, k with $n \geq k$, and let $\Omega = [n]$. How large can $|\mathcal{F}|$ be before $\text{VC}(\mathcal{F}) = k$? If \mathcal{F} consists of all sets of size at most $k - 1$, then it cannot shatter a set of size k . In this case,

$$|\mathcal{F}| = \sum_{i=0}^{k-1} \binom{n}{i}$$

It turns out that this is the largest cardinality possible.

Theorem 4.6 (Sauer-Shelah Lemma). *Suppose $\mathcal{F} \subseteq \mathcal{P}([n])$ has cardinality*

$$|\mathcal{F}| > \sum_{i=0}^{k-1} \binom{n}{i}$$

for some $k \leq n$. Then, \mathcal{F} shatters a subset of $[n]$ of size k .

Proof. We induct on n . If $n = k = 1$, then $\text{VC}(\mathcal{F}) = 1$ if and only if $|\mathcal{F}| = 2$. So suppose $n > 1$ and that the result holds for any smaller sets.

Given a family \mathcal{F} , we create two families \mathcal{F}_1 and \mathcal{F}_2 of sets in $[n - 1]$ as follows:

$$\begin{aligned}\mathcal{F}_1 &:= \{A \subseteq [n - 1] : A \in \mathcal{F} \text{ or } A \cup \{n\} \in \mathcal{F}\} \\ \mathcal{F}_2 &:= \{A \subseteq [n - 1] : A \in \mathcal{F} \text{ and } A \cup \{n\} \in \mathcal{F}\}\end{aligned}$$

Note that the condition in \mathcal{F}_1 is not exclusive, so $\mathcal{F}_2 \subseteq \mathcal{F}_1$.

We claim that

$$|\mathcal{F}| = |\mathcal{F}_1| + |\mathcal{F}_2|$$

Clearly,

$$|\mathcal{F}_1| = \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}_1}, \quad |\mathcal{F}_2| = \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}_2}$$

and

$$|\mathcal{F}| = \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}} + \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}}$$

so it suffices to check that for every $A \subseteq [n - 1]$,

$$\sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}} + \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}} = \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}_1} + \sum_{A \subseteq [n-1]} \mathbf{1}_{A \in \mathcal{F}_2}$$

If only one of A and $A \cup \{n\}$ are in \mathcal{F} , then A is in \mathcal{F}_1 , but not \mathcal{F}_2 . If both A and $A \cup \{n\}$ are in \mathcal{F} , then A is in both \mathcal{F}_1 and \mathcal{F}_2 . In either case, A is counted the same number of times on each side, so the equality holds.

Now, we have

$$\sum_{i=0}^{k-1} \binom{n}{i} = \sum_{i=0}^{k-1} \binom{n-1}{i} + \sum_{i=i}^{k-1} \binom{n-1}{i-1}$$

$$= \sum_{i=0}^{k-1} \binom{n-1}{i} + \sum_{i=0}^{k-2} \binom{n-1}{i}$$

by the recursion formula $\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}$ for binomial coefficients.

So, if $|\mathcal{F}| > \sum_{i=0}^{k-1} \binom{n}{i}$, then either $|\mathcal{F}_1| > \sum_{i=0}^{k-1} \binom{n-1}{i}$ or $|\mathcal{F}_2| > \sum_{i=0}^{k-2} \binom{n-1}{i}$.

In the first case, the inductive hypothesis gives that \mathcal{F}_1 shatters a subset of $[n-1]$ of size k , in which case, \mathcal{F} shatters the same subset.

In the second case, the family \mathcal{F}_2 shatters a subset S of $[n-1]$ of size $k-1$. For each set $B \in \mathcal{F}_2$, both B and $B \cup \{n\}$ are in \mathcal{F} , so \mathcal{F} shatters $S \cup \{n\}$, which has size k . ■

We have a strengthening of the Sauer-Shelah theorem that implies that \mathcal{F} in fact shatters at least $|\mathcal{F}|$ sets.

Theorem 4.7 (Pajor). *Suppose $\mathcal{F} \subseteq \mathcal{P}([n])$ has cardinality*

$$|\mathcal{F}| > \sum_{i=0}^{k-1} \binom{n}{i}$$

for some $k \leq n$. Then, \mathcal{F} shatters at least $|\mathcal{F}|$ sets.

This theorem immediately implies the Sauer-Shelah lemma, since only $\sum_{i=9}^{k-1} \binom{n}{i} < |\mathcal{F}|$ subsets have cardinality less than k .

Proof. We induct on n . If $n = 0$, the sum is empty. But, every family of only one set already shatters the empty set. So, suppose $n > 1$ and that the result holds for any smaller sets.

Given a family \mathcal{F} satisfying the hypotheses of the result for n , we split \mathcal{F} into disjoint subfamilies, \mathcal{F}_1 and \mathcal{F}_2 , where \mathcal{F}_1 contains all the subsets containing n , and \mathcal{F} is its complement, containing all the subsets that do not contain n .

By the inductive hypothesis, \mathcal{F}_1 and \mathcal{F}_2 each shatter two collections of sets whose sizes add to at least $|\mathcal{F}|$.

None of the sets S shattered by either family can contain n , since such sets cannot be shattered by \mathcal{F}_1 , since any subset of S not containing n cannot be created by intersection; nor by \mathcal{F}_2 , since any subset of S containing n cannot be created by intersection.

However, some of the shattered sets S may be shattered by both \mathcal{F}_1 and \mathcal{F}_2 . If S is shattered by only one of \mathcal{F}_1 and \mathcal{F}_2 , then it contributes one to the number of sets shattered by the subfamily and also to the number of sets shattered by \mathcal{F} . Otherwise, if S is shattered by both \mathcal{F}_1 and \mathcal{F}_2 , then both S and $S \cup \{x\}$ are shattered by \mathcal{F} , so S contributes two to the number of shattered sets of the subfamilies and of \mathcal{F} .

Thus, \mathcal{F} shatters at least as many sets as the number of set shattered by \mathcal{F}_1 and \mathcal{F}_2 , which is at least $|\mathcal{F}|$. ■

5 Graph Colouring

A *proper (vertex) colouring* of a graph $G = (V, E)$ is a labelling of the vertex set $c : V \rightarrow [n]$ such that $c(u) \neq c(v)$ whenever $(u, v) \in E$, and the elements of $[n]$ are traditionally called *colours*.

The *chromatic number* $\chi(G)$ of a graph G is the minimum n for which such a labelling of G exists.

Example.

- For any n , $\chi(K_n) = n$, since all n vertices are adjacent to every other vertex.
- For any n , $\chi(C_{2n}) = 2$.
- For any n , $\chi(C_{2n+1}) = 3$.
- A graph G is bipartite if and only if $\chi(G) = 2$.

△

We write

$$\Delta(G) := \max_{v \in V(G)} \deg(v)$$

for the maximum degree of G .

Lemma 5.1. *For any graph G ,*

$$\chi(G) \leq \Delta(G) + 1$$

Proof. Pick any vertex of G , and greedily assign it any colour not present amongst its previously picked neighbours, then repeat. It will always be possible to assign a vertex a valid colour, since each vertex has at most $\Delta(G)$ neighbours that can already be coloured, and there are $\Delta(G) + 1$ colours available. ■

As we will see, this bound for the number of colours needed is rarely sharp.

Lemma 5.2. *Let G be a connected graph that has a vertex x of degree $\deg(x) < \Delta(G)$. Then, $\chi(G) \leq \Delta(G)$.*

Proof. For each vertex in G , determine the length of the shortest path to x . Because G is connected, this distance is well-defined.

Let k be the maximum distance, and for each $0 \leq i \leq k$, define the set

$$V_i := \{v \in V : d(v, x) = i\}$$

Each vertex in V_i is adjacent to a vertex in V_{i-1} via an edge along a shortest path to x .

Now, consider the induced subgraph $G[V_k]$. Because every vertex in V_k has at least one edge to a vertex in V_{k-1} , each vertex has degree at most $\Delta(G) - 1$, so by the previous lemma, $G[V_k]$ is $\Delta(G)$ -colourable.

Now, consider the induced subgraph $G[V_{k-1}]$. Again, every vertex in V_{k-1} has at least one edge to a vertex in V_{k-2} , so each vertex has degree at most $\Delta(G) - 1$, so we can greedily colour $G[V_{k-1}]$ with $\Delta(G)$ colours, taking the colours used for the previous layer into account.

The same argument continues for each V_i with $0 < i \leq k$ until only $V_0 = \{x\}$ remains. By assumption, $\deg(x) < \Delta(G)$, so we have a colour left for x . ■

A graph is k -connected if it requires the deletion of at least k vertices to disconnect it.

Example. Any connected graph is at least 1-connected. △

Example. The path graph P_3 is 1-connected but not 2-connected, as deleting the middle vertex disconnects the graph. △

Example. The cycle graph C_4 is 2-connected. △

Theorem 5.3 (Brooks). *If G is a connected graph which is neither complete nor an odd cycle, then $\chi(G) \leq \Delta(G)$. Otherwise, $\chi(G) = \Delta(G) + 1$.*

Proof. If $G = K_n$, then $\chi(G) = n = \Delta(G) + 1$, and if $G = C_{2n+1}$, then $\chi(G) = 3 = \Delta(G) + 1$.

Also, if $\Delta(G) = 1$, then $G = K_2$, and if $\Delta(G) = 2$, then G is either a cycle or a path, and paths have chromatic number 2 via greedy colouring.

Otherwise, assume that G is neither complete nor an odd cycle, and that $\Delta(G) \geq 3$. We split into three cases:

- (i) G is 1-connected but not 2-connected;
- (ii) G is 2-connected but not 3-connected;
- (iii) G is 3-connected.

- (i) Let v be a vertex whose removal disconnects G into the connected components $G \setminus \{v\} = \bigcup_i G_i$.

Consider the induced subgraphs $G[G_i \cup \{v\}]$. In this induced subgraph, v has degree less than $\Delta(G)$, since it has edges to other connected components, so this induced subgraph can be coloured with $\Delta(G)$ many colours via the previous lemma.

By permuting the colourings in each induced subgraph, we can ensure that v has the same colour in each case, so the union of the colourings gives a proper colouring for G .

- (ii) Let u, v be a pair of vertices whose removal disconnects G into the connected components $G \setminus \{u, v\} = \bigcup_i G_i$.

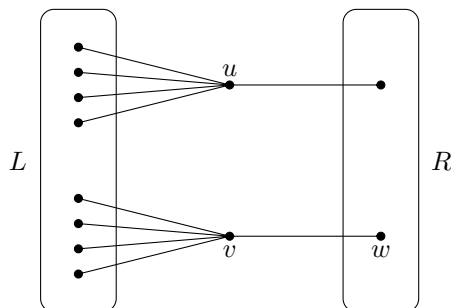
Note that both u and v have at least one edge incident to each component, since if, say, only u has such an edge, then deleting u alone disconnects G , contradicting that G is 2-connected.

Through identical arguments as in case (i), we may colour the induced subgraphs $G[G_i \cup \{u, v\}]$.

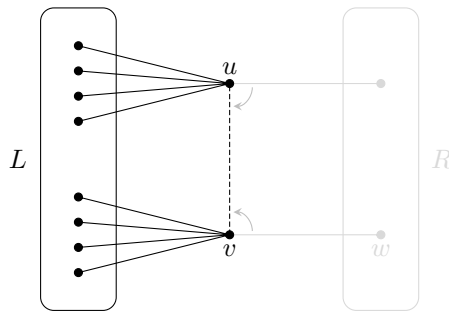
If u and v are adjacent, then in each of the colourings, they are assigned different colours, so by permuting the colourings in each induced subgraph, we can again take the union of these colourings to obtain a proper colouring of G .

Otherwise, u and v are not adjacent. Continue as before, but colour each induced subgraph as though there were an edge connecting u and v . Note that this cannot increase the maximum degree beyond $\Delta(G)$, since u and v previously had at least one other edge to a different connected component.

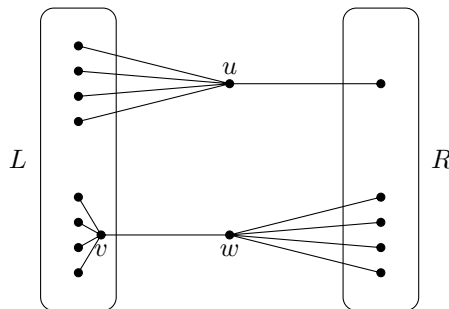
However, this might increase the degree of both u and v to exactly $\Delta(G)$, in which case we may not have a vertex of degree less than $\Delta(G)$ with which to apply the lemma. This happens if and only if $G \setminus \{u, v\} = L \cup R$ has two connected components, and u and v each only have one edge incident to one of them, say R ,



since in this case, adding the edge between u and v in $G[L \cup \{u, v\}]$ leaves their degrees unchanged:



Instead, replace v by its neighbour w in R :



The vertices u and w disconnect G , and now in both components of $G \setminus \{u, w\}$, at least one of u and w has at most $\Delta(G) - 2$ neighbours, so adding in the edge (u, w) leaves each piece with maximum degree $\Delta(G)$, and at least one vertex of smaller degree.

- (iii) We claim that there is an induced path in G of length 2, passing through vertices, say, u, x, v , such that u and v are not adjacent.

Let S be a maximal complete subgraph of G . By assumption, G is not complete, so there is a vertex $u \in G \setminus S$ adjacent to a vertex $x \in S$. There must also be a vertex $v \in S$ not adjacent to u , since if every vertex in S were adjacent to u , then $S \cup \{u\}$ would be a larger complete subgraph, contradicting that S is maximal. This proves the claim.

Colour the vertices u and v with the same colour. Since G is 3-connected, the graph $G \setminus \{u, v\}$ is connected, so each vertex in it has a well-defined distance from x . We proceed as in the proof of the previous lemma, greedily colouring in the subgraphs induced on the sets V_i of vertices at distance i from x . Once we reach x , it has two neighbours u and v with the same colour, so there is a spare colour for x .

■

5.1 The Chromatic Polynomial

Given a graph G , we define the function $P_G : \mathbb{N}_{>0} \rightarrow \mathbb{N}$ as:

$$P_G(k) := \# \text{ of proper } k\text{-colourings of } G$$

where two colourings are considered distinct if there is a vertex labelled with different colours in the two colourings.

Example. For the complete graph K_n on n vertices, we have:

$$P_{K_n}(k) = \prod_{i=0}^{n-1} (k - i)$$

If $k < n$, then K_n is not k -colourable, so $P_{K_n}(k) = 0$. If $k \geq n$, then k -colourings exist: we may select any of the k colours for the first vertex, any of the remaining $k - 1$ for the second, etc.

So, there are $k(k-1)(k-2) \cdots (k-(n-1))$ such colourings, and this formula also agrees with the $k < n$ case since one of the factors would vanish. \triangle

Example. For the path graph P_n on n edges and $n + 1$ vertices, we have:

$$P_{P_n}(k) = k(k-1)^n$$

If $k = 1$ and $n \geq 1$, then there are no colourings. If $k > 1$, we may choose any of the k colours for the first vertex. Then, traversing the path, for each of the n remaining vertices v_{i+1} , we may choose any of the $k - 1$ colours distinct from the colour of the previous vertex v_i .

So, there are $k(k-1)(k-1) \cdots (k-1)$ such colourings, and this formula agrees with the $k = 1$ case, since every factor past the first would vanish. \triangle

Example. For the empty graph E_n on n vertices, we have:

$$P_{E_n}(k) = k^n$$

Since there are no adjacencies, every vertex can be independently coloured with any of the k colours, and there are n vertices, so there are k^n total colourings. \triangle

So far, we have seen that

$$\begin{aligned} P_{K_n}(k) &= \prod_{i=0}^{n-1} (k-i) \\ P_{P_n}(k) &= k(k-1)^n \\ P_{E_n}(k) &= k^n \end{aligned}$$

In each case, P_G is a polynomial in k . This turns out to be true for all finite graphs, and we call P_G the *chromatic polynomial*. This fact is not obvious:

Example. For the cycle graph C_n on n vertices, we can start similarly to the path graph: we choose a vertex v_0 , and colour it with any of the k colours. Then, traversing the cycle in one direction, we can colour each vertex v_{i+1} with any of the $k - 1$ colours distinct from the colour of the previous vertex v_i .

However, how can we colour the vertex v_{n-1} that is adjacent to v_0 ? The number of valid colours depends on whether v_{n-2} is the same colour as v_0 or not. At this point, it is unclear as to how we should proceed. \triangle

Often, when proving a result on all finite graphs, we induct on the size of the vertex set. However, in the inductive step, the number of choices for the new vertex depends not only on the number of neighbours, but also on the colouring on the neighbours.

Instead, we might try to induct on the number of edges. Take a finite graph G and let $e = (x, y)$ be an edge in G . From G , we construct the graph $G \setminus e$ by deleting e , and the graph G/e by contracting e .

We can partition the possible k -colourings c of $G \setminus e$ into two cases:

- $c(x) \neq c(y)$;
- $c(x) = c(y)$.

In the former case, each such colouring is also an admissible colouring of G , since x and y are adjacent in G , but are assigned different colours. In the latter case, these colourings are not admissible. However, such colourings correspond precisely to the proper colourings of G/e , since x and y are the same vertex in the contraction.

This observation provides us with our inductive step.

Theorem 5.4. *For every finite graph $G = (V, E)$ containing an edge $e = (x, y) \in E$, and for every $k \geq 1$,*

$$P_G(k) = P_{G \setminus e}(k) - P_{G/e}(k)$$

Consequently, P_G coincides with a polynomial in $\mathbb{R}[k]$.

Proof. Every k -colouring of $G \setminus e$ either assigns different colours to x and y , in which case, it corresponds to a proper colouring of G , or it assigns them the same colour. So, it suffices to check that the number of colourings of $G \setminus e$ in which x and y are assigned the same colour is equal to the number of colourings of G/e .

Given a colouring of $G \setminus e$ with $c(x) = c(y)$ are the same colour, then we can construct a corresponding colouring of G/e by colouring the contracted vertex as $c(x) = c(y)$, and leaving all other vertices unchanged; conversely, given a colouring of G/e , we can construct a colouring of $G \setminus e$ by colouring x and y the same as the contracted vertex.

We deduce that P_G is a polynomial by induction on $|E|$.

If $|E| = 0$, then $G = E_{|V|}$ is an empty graph, and we have that $P_{E_n}(k) = k^n$ is a polynomial. Otherwise, G has an edge $e = (x, y)$, and $G \setminus e$ and G/e are graphs with fewer edges than G , so $P_G(k) = P_{G \setminus e}(k) - P_{G/e}(k)$ is the difference of two polynomials and is hence a polynomial. ■

We can use this recurrence relation to compute the chromatic polynomial:

Theorem 5.5. *The chromatic polynomial of the cycle graph C_n is*

$$P_{C_n}(k) = (k-1)^n + (-1)^n(k-1)$$

Proof. If $n = 3$, then $C_3 = K_3$, and

$$\begin{aligned} (k-1)^3 + (-1)^3(k-1) &= (k-1)^3 - (k-1) \\ &= k^3 - 3k^2 + 3k - 1 - k + 1 \\ &= k^3 - 3k^2 + 2k \\ &= k(k-1)(k-2) \\ &= P_{K_3}(k) \\ &= P_{C_3}(k) \end{aligned}$$

as required.

Now suppose $n > 3$, and that the result holds for all smaller cycles. Let $e \in E(C_n)$. Then, $C_n \setminus e = P_{n-1}$ and $C_n/e = C_{n-1}$, so

$$\begin{aligned} P_{C_n}(k) &= P_{C_n \setminus e}(k) - P_{C_n/e}(k) \\ &= P_{P_{n-1}}(k) - P_{C_{n-1}}(k) \\ &= [k(k-1)^{n-1}] - [(k-1)^{n-1} + (-1)^{n-1}(k-1)] \\ &= (k-1)(k-1)^{n-1} + (-1)(-1)^{n-1}(k-1) \\ &= (k-1)^n + (-1)^n(k-1) \end{aligned}$$

completing the inductive step. ■

Note that if $k = 2$, then,

$$\begin{aligned} P_{C_n}(2) &= (2-1)^n + (-1)^n(2-1) \\ &= 1 + (-1)^n \end{aligned}$$

$$= \begin{cases} 0 & n \text{ odd} \\ 2 & n \text{ even} \end{cases}$$

There are some other properties of the chromatic polynomial that we can immediately deduce from the recurrence relation.

Theorem 5.6. *For any finite graph $G = (V, E)$,*

- (i) *The degree of P_G is $|V|$;*
- (ii) *P_G is monic;*
- (iii) *The coefficients of P_G have alternating signs;*
- (iv) *$P_G(0) = 0$.*

Proof. In all cases, we induct on $|E|$ to use the chromatic polynomial recurrence relation.

If $|E| = 0$, then G is empty and $P_{E|V|}(k) = k^{|V|}$ is a monic polynomial of degree $|V|$. Also, the coefficients are all zero apart from this first term, so they trivially alternate signs. We also have $P_{E|V|}(0) = 0^{|V|} = 0$. This establishes the base case for all four claims.

Now, suppose $|E| > 0$, and that the result holds for all graphs with fewer edges, and let $e \in E$, so

$$P_G(k) = P_{G \setminus e}(k) - P_{G/e}(k)$$

- (i),(ii) The graph $G \setminus e$ has fewer edges than G , but the same number of vertices, so $P_{G \setminus e}(k)$ is a monic polynomial of degree $|V(G \setminus e)| = |V|$ by the inductive hypothesis. The graph G/e has fewer edges and fewer vertices than G , so its chromatic polynomial does not contribute to the $|V|$ th order term in P_G . So, P_G is a monic polynomial of degree $|V|$.
- (iii) The graphs $G \setminus e$ and G/e have fewer edges than G , so their chromatic polynomial coefficients have alternating signs by the induction hypothesis. $P_{G/e}$ also has degree $|V| - 1$ by property (i), so its signs are opposite to that of $P_{G \setminus e}$, and this alternation is preserved in their difference.
- (iv) By the inductive hypothesis, $P_{G \setminus e}(0) = P_{G/e}(0) = 0$. So $P_G(0) = 0 - 0 = 0$. ■

Thus, if $G = (V, E)$ with $|V| = n$, then its chromatic polynomial is of the form:

$$P_G(k) = k^n - c_{n-1}k^{n-1} + c_{n-2}k^{n-2} - c_{n-3}k^{n-3} + \dots$$

Theorem 5.7. *Let $G = (V, E)$ be a finite graph with chromatic polynomial*

$$P_G(k) = k^n - c_{n-1}k^{n-1} + c_{n-2}k^{n-2} - c_{n-3}k^{n-3} + \dots$$

Then,

(i)

$$c_{n-1} = |E|$$

(ii) and

$$c_{n-2} = \binom{|E|}{2} - T(G)$$

where $T(G)$ is the number of triangles in G .

Proof. We induct on $|E|$. For $|E| = 0$, $P_{E_n}(k) = k^n$, and the two coefficients are 0, as required. Now, suppose $|E| > 0$, and that the result holds for all graphs with fewer edges.

Let $e = (x, y) \in E$ and let $c_n(G)$ denote the coefficient the n th degree term of P_G .

(i) Comparing the $(n - 1)$ th degree terms of the three polynomials, we have

$$[-c_{n-1}(G)] = [-c_{n-1}(G \setminus e)] - c_{n-1}(G/e)$$

Note that the first two terms are negative because c_{n-1} is the *second* coefficient of P_G and $P_{G \setminus e}$, but the *leading* coefficient of $P_{G/e}$, since $P_{G/e}$ degree one less than P_G .

By the previous theorem, the leading coefficient $c_{n-1}(G/e)$ is 1, and by the inductive hypothesis, $c_{n-1}(G \setminus e) = |E(G \setminus e)| = |E| - 1$. So,

$$\begin{aligned} c_{n-1}(G) &= (|E| - 1) + 1 \\ &= |E| \end{aligned}$$

This completes the inductive step.

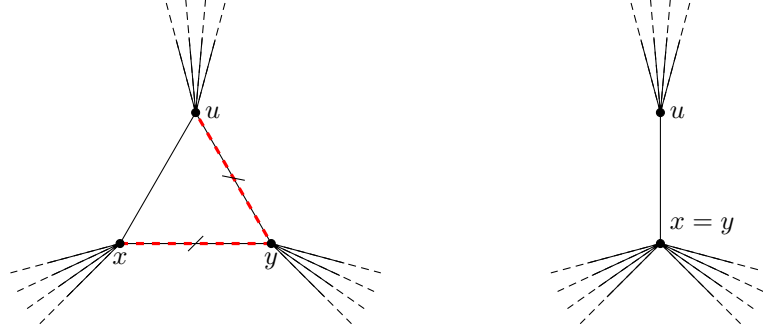
(ii) Comparing the $(n - 2)$ th degree terms of the three polynomials, we have

$$c_{n-2}(G) = c_{n-2}(G \setminus e) - [-c_{n-2}(G/e)]$$

This time, c_{n-2} is the third coefficient of P_G and $P_{G \setminus e}$, which is positive, but the second coefficient of $P_{G/e}$, which is negative. By part (i),

$$c_{n-2}(G/e) = |E(G/e)|$$

When we contract $e = (x, y)$, we lose one edge from E , but we also lose an edge for every vertex u adjacent to both x and y , since those two edges are combined into a single edge in the contraction:



Let T_0 be the number of such vertices, so

$$c_{n-2}(G/e) = |E| - 1 - T_0$$

Then, by the inductive hypothesis,

$$\begin{aligned} c_{n-2}(G \setminus e) &= \binom{|E(G \setminus e)|}{2} - T(G \setminus e) \\ &= \binom{|E| - 1}{2} - T(G \setminus e) \end{aligned}$$

$T(G \setminus e)$ is equal to $T(G)$ minus the number of triangles that were removed when deleting e . That is, the number of triangles that contain e . But this is exactly T_0 , so

$$c_{n-2}(G \setminus e) = \binom{|E| - 1}{2} - (T(G) - T_0)$$

So the c_{n-2} coefficient of G is:

$$\begin{aligned}
 c_{n-2}(G) &= c_{n-2}(G \setminus e) + c_{n-2}(G/e) \\
 &= \binom{|E| - 1}{2} - (T(G) - T_0) + |E| - 1 - T_0 \\
 &= \binom{|E| - 1}{2} + |E| - 1 - T(G) \\
 &= \binom{|E|}{2} - T(G)
 \end{aligned}$$

as required. ■

6 Matroids

A *matroid* (E, \mathcal{I}) consists of a *ground set* E , and a family $\mathcal{I} \subseteq \mathcal{P}(E)$ of its subsets satisfying

- (i) $\emptyset \in \mathcal{I}$;
- (ii) If $A \subseteq B$ and $B \in \mathcal{I}$, then $A \in \mathcal{I}$ (*hereditary property* or *downward-closedness*);
- (iii) If $A, B \in \mathcal{I}$ and $|A| > |B|$, then there is an element $a \in A$ such that $\{a\} \cup B \in \mathcal{I}$ (*exchange condition*).

The sets in \mathcal{I} are called the *independent sets* of the matroid.

Example. Let V be a vector space and $E \subseteq V$ be a set of vectors. If \mathcal{I} is the collection of linearly independent subsets of E , then (E, \mathcal{I}) is a matroid called the *vector matroid*. △

Example. Fix integers n, r with $n < r$. Let $E = [n]$ and take $\mathcal{I} = \{S \subseteq E : |S| < r\}$ to be the set of subsets with cardinality at most r . Then, (E, \mathcal{I}) is a matroid called the *uniform matroid* $U_{n,r}$. △

Lemma 6.1 (Characterisation of Trees). *For any graph $G = (V, E)$, any two of the following imply the third (and hence that G is a tree):*

- (i) $|E| = |V| - 1$;
- (ii) G is connected;
- (iii) G is acyclic.

Theorem 6.2 (Graphic Matroids). *Let $G = (V, E)$ be a graph, and let $\mathcal{I} \subseteq \mathcal{P}(E)$ be the set of acyclic subsets of E (i.e. the set of forests). Then, (E, \mathcal{I}) is a matroid.*

Proof. The empty set is acyclic, and any subset of an acyclic subset is still acyclic. Now, let $A, B \in \mathcal{I}$ with $|A| > |B|$, and consider the graph $G_B = (V, B)$. B is acyclic, so G_B is a forest.

If there is an edge $a \in A$ that connects distinct components of B , then we are done, as $\{a\} \cup B$ is acyclic. Otherwise, suppose there are no such edges. But then, every edge in A lies within the connected components of G_B . Since A is acyclic, it cannot have more edges in each component than a tree does, so it cannot have more elements than B . ■

A matroid (E, \mathcal{I}) is *representable* over a field K if there is a vector space V over K and a map $\phi : E \rightarrow V$ such that for each $A \subseteq E$, $A \in \mathcal{I}$ if and only if $\phi(A)$ is a linearly independent set in V .

Example. The uniform matroid $U_{4,2}$ cannot be represented over \mathbb{Z}_2 .

Suppose otherwise, so there are 4 vectors x_1, \dots, x_4 such that any pair are linearly independent, but any three are linearly dependent. The only possible linear dependency of three vectors is of the form

$$x_1 + x_2 + x_3 = 0$$

since the only coefficients available are 0 and 1..

But then, adding the linear dependencies $x_1 + x_2 + x_3 = 0$ and $x_1 + x_2 + x_4 = 0$ yields $x_3 + x_4 = 0$, contradicting the linear independence of any pair of vectors. \triangle

Example. Graphic matroids can be represented over any field. \triangle

Example. $U_{n,r}$ is representable over \mathbb{R} for any n and r . \triangle

6.1 Rado's Theorem

We recall the set-theoretic statement of Hall's theorem.

A *transversal* or *system of distinct representatives (SDR)* of a family of subsets $\mathcal{F} \subseteq \mathcal{P}(X)$ is a subset of X obtained by selecting a distinct representative from each subset $S \in \mathcal{F}$.

Theorem 6.3 (Hall). *Let $\{S_i\}_{i=1}^n$ be a collection of subsets of a set X . Then, there is a transversal of $\{S_i\}$ if and only for every subset of indices $\sigma \subseteq [n]$, we have*

$$\left| \bigcup_{i \in \sigma} S_i \right| \geq |\sigma|$$

Proof. Apply the graph-theoretic variant of Hall's theorem on the bipartite graph $G = (L \cup R, E)$, where $L = \{S_i\}$ and $R = X$, and $(S_i, e) \in E$ if and only if $e \in S_i$. \blacksquare

Suppose we have sets $S_i \subseteq E$ in a matroid (E, \mathcal{I}) . Under what conditions can we find a transversal of the S_i that is an independent set?

For a set $A \subseteq E$, we define the *rank* $r(A)$ of A to be the cardinality of the largest independent set in A :

$$r(A) := \max\{|B| : B \subseteq A, B \in \mathcal{I}\}$$

Example. If (E, \mathcal{I}) is a vector matroid, then $r(E) = \dim(E)$, and $r(A) = \dim(\text{span}(A))$. \triangle

Lemma 6.4 (Rank Submodularity). *Given a matroid (E, \mathcal{I}) , the rank function satisfies, for any $A, B \subseteq E$:*

$$r(A \cup B) + r(A \cap B) \leq r(A) + r(B)$$

Proof. Choose a maximal independent set $I_\cap \in \mathcal{I}$ in the intersection $A \cap B$. By definition of matroid rank, $|I_\cap| = r(A \cap B) =: n$. Through repeated applications of the exchange condition, extend this set to a maximal independent set I_\cup in $A \cup B$ of size $|I_\cup| = r(A \cup B) =: m$.

Let $a := |I_\cup \setminus B|$ and $b := |I_\cup \setminus A|$, so

$$\begin{aligned} |I_\cup| &= |I_\cup \setminus B| + |I_\cup \setminus A| + |I_\cup \cap (A \cap B)| \\ |I_\cup| &= |I_\cup \setminus B| + |I_\cup \setminus A| + |I_\cap| \\ m &= a + b + n \end{aligned}$$

Then,

$$r(A \cup B) + r(A \cap B) = m + n$$

$$\begin{aligned}
&= (a + b + n) + n \\
&= (a + n) + (b + n)
\end{aligned}$$

Note that $I_\cap \sqcup (I \cap \setminus B) = (I_\cup \cap A) \subseteq A$ is an independent subset of A of size $a + n$. The rank $r(A)$ is defined to be the size of a maximal independent subset of A , so we have $a + n \leq r(A)$. Similarly, $b + n \leq r(B)$, giving:

$$\leq r(A) + r(B)$$

as required. ■

Theorem 6.5 (Rado). *Let (E, \mathcal{I}) be a matroid, and let $S_1, \dots, S_n \in \mathcal{P}(E)$ be arbitrary subsets of E . If for every set $\sigma \subseteq [n]$ of indices,*

$$r\left(\bigcup_{i \in \sigma} S_i\right) \geq |\sigma|$$

then there is an independent transversal. That is, a set $\{e_1, \dots, e_n\} \in \mathcal{I}$ of n distinct elements of E with $e_i \in S_i$ for each i .

Theorem 6.6 (Horn). *Let $X = \{x_1, \dots, x_n\}$ be vectors in a vector space V , and suppose that for each set $\sigma \subseteq [n]$ of indices,*

$$\dim(\text{span}(\{x_i : i \in \sigma\})) \geq \frac{|\sigma|}{2}$$

Then, the set of vectors can be partitioned into two linearly independent sets.

Proof. Let $E = X \sqcup X$, and denote the second copies of the x_i by x'_i . We declare a subset of E as independent in \mathcal{I} if its x_i elements are linearly independent and its x'_i elements are linearly independent.

The empty set is independent in both cases, and a subset of a linearly independent set is still linearly independent, so (E, \mathcal{I}) satisfies downward-closure. Then, if $A, B \in \mathcal{I}$ and $|A| > |B|$, then A has more x_i than B , or A has more x'_i than B (or both). By considering only the larger set, this is effectively the ordinary vector matroid, so the exchange condition holds similarly in any case.

For each $i \in [n]$, define the set $S_i = \{x_i, x'_i\}$. Let $\sigma \subseteq [n]$. By the hypotheses of the theorem, there is a subset $\tau \subseteq \sigma$ of at least half the size for which the vectors $\{x_i : i \in \tau\}$ are linearly independent. But then,

$$\bigcup_{i \in \sigma} S_i \supseteq \bigcup_{i \in \tau} S_i$$

which has at least $|\sigma|$ elements, and is independent in the matroid. So, the sets S_i satisfy the hypotheses of Rado's theorem, so there is an independent transversal of the S_i . That is, an independent selection of exactly one x_i or x'_i from each S_i . Then, the set of selected x_i and the set of selected x'_i give the required partition. ■

7 Random Graphs

Given $r \in \mathbb{N}$, we define $R(r, r)$ to be the smallest positive integer such that for every edge 2-colouring of K_n , there is a monochromatic K_r as a subgraph.

Theorem 7.1. $R(3, 3) = 6$.

Proof. Suppose the edges of K_6 are coloured red and blue. Select a vertex u . There are five edges incident to u , so by the pigeonhole principle, at least three of these edges (u, v_1) , (u, v_2) , and (u, v_3) are the same colour, say, red. If any of the edges connecting the v_i are red, then this forms a red triangle including u . Otherwise, none of the edges are red, in which case, the v_i form a blue monochromatic triangle. ■

Theorem 7.2 (Erdős Lower Bound for $R(r, r)$). *Let $r \geq 3$. Then,*

$$R(r, r) \geq 2^{\frac{r-1}{2}}$$

Proof. Colour the edges of K_n red or blue independently with probability $1/2$ each. For any fixed set S_i of r vertices, define the random variable $X(S_i)$ to be 1 if the K_r induced on S_i is monochromatic, and 0 otherwise. For any S_i , the expectation of $X(S_i)$ is the probability that all $\binom{r}{2}$ edges are the same colour:

$$\mathbb{E}[X(S_i)] = 2 \cdot \left(\frac{1}{2}\right)^{\binom{r}{2}} = 2^{1-\binom{r}{2}}$$

There are $\binom{n}{r}$ many possible subsets S_i , so the number of monochromatic K_r is the sum

$$\sum_{i=1}^{\binom{n}{r}} X(S_i)$$

which has expected value

$$\begin{aligned} \mathbb{E} \left[\sum_{i=1}^{\binom{n}{r}} X(S_i) \right] &= \sum_{i=1}^{\binom{n}{r}} \mathbb{E}[X(S_i)] \\ &= \binom{n}{r} \cdot 2^{1-\binom{r}{2}} \\ &= \frac{2 \cdot n!}{r!(n-r)!2^{\frac{r(r-1)}{2}}} \\ &< \frac{2 \cdot n^r}{r!2^{\frac{r(r-1)}{2}}} \end{aligned}$$

If this is less than 1, then there are colourings without any monochromatic K_r . But, if $n \leq 2^{\frac{r-1}{2}}$, then this expectation is at least 1. ■

Given $n \in \mathbb{N}$ and $p \in [0, 1]$, we write $G \sim G(n, p)$, or just $G_{n, p}$, if G is a random graph with vertex set $E(G) = [n]$ and each possible edge is included independently at random in $E(G)$ with probability p .

Lemma 7.3. *For any $p \in (0, 1)$ and any m ,*

$$\exp\left(\frac{-mp}{1-p}\right) \leq (1-p)^m \leq \exp(-mp)$$

Lemma 7.4. *For $1 \leq k \leq n$,*

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

Theorem 7.5 (Linearity of Expectation). *For any random variables X_1, \dots, X_n , which may be dependent, and constants c_1, \dots, c_n ,*

$$\mathbb{E} \left[\sum_{i=1}^n c_i X_i \right] = \sum_{i=1}^n c_i \cdot \mathbb{E}[X_i]$$

Example. We compute the expected number of triangles in $G_{n,p}$ using the linearity of expectation. There are $\binom{n}{3}$ possible triangles, and each triangle has probability p^3 of being included. For each triangle T , define the random variable X_T to be 1 if T is in $G_{n,p}$, and 0 otherwise. Then, the number of triangles in $G_{n,p}$ is

$$\sum_{T \in G} X_T$$

The X_T are not independent, since if T is in G , then any triangle sharing an edge with T is more likely to also be in G . But by the linearity of expectation, the expected number of triangles in G is

$$\begin{aligned} \mathbb{E} \left[\sum_{T \in G} X_T \right] &= \sum_{T \in G} \mathbb{E}[X_T] \\ &= \binom{n}{3} p^3 \end{aligned}$$

△

Theorem 7.6 (Markov's Inequality). *If X is a non-negative random variable, and $t > 0$, then*

$$\mathbb{P}(X > t) \leq \frac{\mathbb{E}[X]}{t}$$

7.1 Chromatic Numbers

Theorem 7.7 (Chromatic Number of a Random Graph). *Let $k \in \mathbb{Z}$ and suppose p satisfies*

$$p \geq \frac{2k \log(k) + 4k}{n}$$

Then,

$$\mathbb{P}(\chi(G_{n,p}) \leq k) \leq \exp\left(-\frac{n}{2k}\right)$$

If a graph is k -coloured, then one of the colour classes has at least $\frac{n}{k}$ vertices in it by the pigeonhole principle, so it suffices to prove that if r is an integer satisfying

$$\frac{n}{k} \leq r \leq \frac{n}{k} + 1$$

then the probability that $G_{n,p}$ contains an independent set of r vertices is less than $\exp(-r/2)$.

Theorem 7.8 (Independence Number of a Random Graph). *Let $k \in \mathbb{Z}$ and suppose p satisfies*

$$p \geq \frac{2k \log(k) + 4k}{n}$$

Then,

$$\mathbb{P}(G_{n,p} \text{ has an independent set of size at least } r) \leq \exp\left(-\frac{r}{2}\right)$$

Lemma 7.9. *Let $g, n \in \mathbb{Z}$ and $p \in [0,1]$ satisfy*

$$\frac{5}{n} \leq p \leq \frac{n^{\frac{1}{g}}}{n}$$

and let X be the number of cycles of length at most $g-1$ in $G_{n,p}$. Then, $\mathbb{E}(X) \leq \frac{n}{4}$.

Corollary 7.9.1. *For large n , there is a graph on at most n vertices with chromatic number at least*

$$\frac{\log(n)}{4 \log(\log(n))} - 1$$

and no cycles shorter than

$$\frac{\log(n)}{\log(\log(n))} - 1$$

7.2 Connectedness

Theorem 7.10 (Connectedness of Random Graphs). *Let (c_n) be a sequence, and let*

$$p(n) = \frac{\log(n)}{n} + \frac{c_n}{n}$$

Then,

$$\mathbb{P}(G_{n,p} \text{ is connected}) \rightarrow \begin{cases} 0 & c_n \rightarrow -\infty \\ 1 & c_n \rightarrow \infty \end{cases}$$