



Tel Aviv University  
The Blavatnik School of Computer Science  
Winter 2018

# ABSTRACT ALGEBRA IN THEORETICAL COMPUTER SCIENCE

Gil Cohen

# Preface

# Table of Contents

## Lecture 1: Why Abstract Algebra?

1.1	Error correcting codes . . . . .	1-2
1.1.1	Definitions . . . . .	1-2
1.1.2	Lower and Upper Bounds . . . . .	1-3
1.1.2.1	The Reed-Solomon Error Correcting Code . . . . .	1-3
1.1.2.2	Remarks . . . . .	1-4
1.2	Expander Graphs . . . . .	1-5
1.2.1	The Construction . . . . .	1-6
1.3	Private Information Retrieval (PIR) . . . . .	1-6
1.3.1	Using Two Servers . . . . .	1-8
1.4	Randomness Mergers . . . . .	1-8
1.4.1	The setup . . . . .	1-8
1.4.2	What cannot be achieved . . . . .	1-10
1.4.3	What can be achieved . . . . .	1-10

---

## Lecture 2: The Fundamental Theorem of Algebra; Groups - The Basics

2.1	The Familiar Number Systems . . . . .	2-1
2.1.1	Enter zero . . . . .	2-1
2.1.2	Negative numbers . . . . .	2-2
2.1.3	Be rational . . . . .	2-2
2.1.4	Get real . . . . .	2-3
2.1.5	Complex numbers . . . . .	2-3
2.2	Will this ever end? . . . . .	2-5
2.3	What else is cool about $\mathbb{C}$ ? . . . . .	2-6
2.4	What does a Turing Machine think of $\mathbb{C}$ ? . . . . .	2-7
2.5	Bézout's Theorem . . . . .	2-7
2.6	Quick Introduction to Group Theory . . . . .	2-8
2.7	The definition of a Group . . . . .	2-9
2.7.1	The associative law . . . . .	2-9
2.7.2	Neutral element . . . . .	2-9
2.7.3	Inverses . . . . .	2-10
2.7.4	The formal definition . . . . .	2-10
2.7.5	But why this definition? . . . . .	2-11

2.8	Commutative groups . . . . .	2-11
2.9	Constructing some groups . . . . .	2-12
2.10	Some basic general properties . . . . .	2-13
2.11	Subgroups . . . . .	2-14
2.12	Cosets . . . . .	2-15
2.13	Normal groups and quotient groups . . . . .	2-16

---

### Lecture 3: Groups - The Basics

3.1	The definition of a Group . . . . .	3-1
3.1.1	The associative law . . . . .	3-2
3.1.2	Neutral element . . . . .	3-2
3.1.3	Inverses . . . . .	3-2
3.1.4	The formal definition . . . . .	3-3
3.1.5	But why this definition? . . . . .	3-3
3.2	Commutative groups . . . . .	3-4
3.3	Constructing some groups . . . . .	3-4
3.4	Some basic general properties . . . . .	3-6
3.5	Subgroups . . . . .	3-6
3.6	Cyclic groups and generated subgroups . . . . .	3-7
3.7	Multiplicative Group . . . . .	3-10
3.8	Normal subgroups . . . . .	3-12
3.9	Homomorphisms . . . . .	3-13

---

### Lecture 4: Quotient Group

4.1	Group Theory - Continued . . . . .	4-1
4.1.1	From previous lecture . . . . .	4-1
4.1.2	First Isomorphism Theorem . . . . .	4-1
4.2	Ring Theory . . . . .	4-2
4.3	Properties of Commutative Rings . . . . .	4-4
4.3.1	Gaussian Integers . . . . .	4-6
4.3.2	The orthogonality of addition and multiplication . . . . .	4-7

---

## Lecture 5: Group Homomorphism

5.1	Fields and PID . . . . .	5-1
5.2	Primes and Irreducibles . . . . .	5-1
5.3	Constructing ideals from existing ones . . . . .	5-2
5.4	Maximal Ideals and Prime Ideals . . . . .	5-4
5.5	Gaussian Integers . . . . .	5-5

---

## Lecture 6: Fields

6.1	Ideal Properties - Continued . . . . .	6-1
6.1.1	From Previous Lecture . . . . .	6-1
6.1.2	Gaussian Integers . . . . .	6-1
6.1.3	$\mathbb{Z}[\sqrt{-5}]$ . . . . .	6-2
6.2	Fields . . . . .	6-4
6.2.1	Fields of Fractions . . . . .	6-4
6.2.2	The Polynomial Ring . . . . .	6-5
6.2.3	Building Finite Fields . . . . .	6-6
6.2.3.1	Finite field of 4 . . . . .	6-6
6.2.3.2	Finite field of 8 . . . . .	6-7
6.2.3.3	Finite field of 9 . . . . .	6-8
6.2.4	Building The Complex Field . . . . .	6-8

---

## Lecture 7: Field Extensions

7.1	Recap - Constructing $\mathbb{C}$ from $\mathbb{R}$ . . . . .	7-1
7.2	Field Extentions . . . . .	7-2

---

## Lecture 8: Fields and Polynomials

8.1	Proving final locations for a winning Peg Solitaire game . . . . .	8-1
8.2	Back to Fields . . . . .	8-4
8.3	Polynomials . . . . .	8-6
8.3.1	Roots of polynomials . . . . .	8-9
8.3.2	Derivatives: formal derivative . . . . .	8-10

---

## Lecture 9: Finite Fields cont.; Small-Biased Sets

9.1	Extension Fields	9-1
9.1.1	Recap	9-1
9.1.2	Splitting Fields	9-1
9.1.3	$F(a_1, \dots, a_n)$	9-2
9.2	Finite Fields	9-2
9.2.1	Existence of $\mathbb{F}_{p^n}$	9-2
9.2.2	Construction of $\mathbb{F}_{p^n}$	9-3
9.3	Small Biased Sets	9-4
9.3.1	Remarkable results	9-5
9.3.2	The powery construction (AGHP)	9-5

---

## Lecture 10: Small bias sets

10.1	Bezout Theorem	10-1
10.2	Ben-Aroya, Ta-Shma Construction	10-1
10.2.1	The construction	10-1
10.2.2	Analyzing the Construction	10-2
10.2.3	Missing Proofs	10-4

---

## Lecture 11: Randomness Mergers

11.1	Small-Biased Sets	11-1
11.1.1	Introduction	11-1
11.1.2	The Probabilistic Method	11-1
11.1.3	Existence proof for a small-biased set	11-2
11.2	Mergers	11-3
11.2.1	Motivation	11-3
11.2.2	Close to Uniform	11-4
11.2.3	Definition	11-5
11.3	Constructing a Merger	11-5
11.3.1	Merging 2 random variables	11-6
11.3.1.1	Schwartz-Zippel lemma	11-7
11.3.1.2	The polynomial method	11-8
11.3.1.3	$Merg$ is $(\rho, \epsilon)$ -random	11-9
11.3.1.4	Fixing Parameters	11-11
11.3.2	Merging multiple random variables	11-12
11.3.2.1	Iterative pair merging	11-12

11.3.2.2 Curve Merger . . . . .	11-12
---------------------------------	-------

---

## Lecture 12: Two Source Extractors and Unbalanced Expanders

12.1 Two Source Extractors . . . . .	12-1
12.1.1 Definitions . . . . .	12-1
12.1.2 Main Results And Constructions . . . . .	12-2
12.2 Unbalanced Expanders . . . . .	12-3
12.2.1 Introduction . . . . .	12-3
12.2.2 Proof Of Existence . . . . .	12-4
12.2.3 An Explicit Polynomial Based Construction . . . . .	12-6

# LECTURE 1

## WHY ABSTRACT ALGEBRA?

---

It is often the case that natural algorithmic questions lead to natural combinatorial questions. A priori it didn't have to be the case, but it is. The reason is unclear to me—whether it is due to the mathematical oriented researchers in computer science or perhaps there is a deeper reason anchored at the problems themselves.

Philosophical aspects aside, in many cases, the algorithmic questions ask for an explicit construction of a certain combinatorial object such as a graph with some desired properties. By explicit here we typically mean that the object of interest can be efficiently generated, that is, there is an algorithm that given, say, the size of the object (i.e., the number of vertices of a graph), outputs the object in polynomial-time (with respect to its size).

In many cases, explicit constructions of combinatorial objects make use of objects and tools from other branches of mathematics, not just combinatorics. The most relevant to computer science seems to be algebra, in particular, linear algebra and basic algebraic structures such as fields, rings, groups, and polynomials over these structures. I.e., abstract algebra. If one would erase every passage that involves algebra from the computer science literature, I'm not sure you would be able to download this file. What I am trying to say, if you haven't noticed, is that algebraic structures have a key role in theoretical computer science. Typically, they serve as tools for constructing and analyzing combinatorial objects that, in turn, are used in algorithms. Sometimes, however, the algebra “appears” only in the analysis (though, of course, it guides the researcher who constructs the object). The best way to convince you of that is to see this in action but we will have to develop the mathematical theory before we can do that.

In this lecture, we will present four examples where algebraic structures come out of nowhere to save the day. We will show the problems. The solutions will have to wait until we develop the algebraic theory we need. This is going to be a recurring theme in the course. We will see problems from different branches of theoretical computer science, coding theory, cryptography, etc which will allow us to get a taste of what is done in these areas. Then, we will see how algebraic structures and reasoning can be used to solve these problems. At first, it will seem to come out of nowhere. As you will see more and more examples, it will seem more natural. If you see enough examples, however, the whole thing (to me at least) starts to look suspicious again.



## 1.1 Error correcting codes

What is an error correcting code? Consider the following scenario: Alice wants to send a message to Bob. The message is sent over some channel, that is not perfect: when sending a message over it, some of the received bits at the other end might be wrong (aka flipped).

**The Coding Problem** asks: what kind of redundancy Alice should add to the message, so that Bob will still be able to extract the message?

We would like to find an encoding function  $C : \{0, 1\}^k \rightarrow \{0, 1\}^n$ , such that given a message  $m \in \{0, 1\}^k$ , if  $C(m) \in \{0, 1\}^n$  is sent over a noisy channel and at most  $(0.1)n$  of the received bits are flipped, the original message  $m$  is still unambiguously recoverable from it.

Clearly,  $C$  must be injective (thus implying, among other things  $n \geq k$ ), but this is obviously not enough, since we want to map different messages  $m_1 \neq m_2$  to strings in  $\{0, 1\}^n$  that are **far apart**. On the other hand, if our code encodes every message in  $\{0, 1\}^k$  as a string in, say,  $\{0, 1\}^{2^k}$  then our "overhead" would be huge - we would have to transmit exponentially many more bits than our original message. Our goal then is to find such an encoding function where the tradeoff between the redundancy of bits transmitted and the distance between encoded words is optimal.

We now give some formal definitions for the task at hand.

### 1.1.1 Definitions

**Definition 1.1.** (Hamming distance) Let  $x, y \in \Sigma^\ell$  be two words of length  $\ell$  over some alphabet  $\Sigma$ , then we define the **Hamming Distance** between the words as

$$d(x, y) \stackrel{\text{def}}{=} |\{i : x_i \neq y_i\}|$$

**Definition 1.2.** (Code) A function  $C : \Sigma_{in}^k \rightarrow \Sigma_{out}^n$  is called a Code with relative distance  $\delta$  and relative rate  $\rho$  if:

- Distance: For any  $m_1 \neq m_2 \in \{0, 1\}^k$ ,

$$d(C(m_1), C(m_2)) \geq \delta \cdot n$$

- Rate: It holds that

$$\frac{k \cdot \log |\Sigma_{in}|}{n \cdot \log |\Sigma_{out}|} \geq \rho$$

In the case where  $\Sigma_{out}^n$  is a vector space (e.g. -  $\{0, 1\}^n$ ) we say that  $C$  is a **linear code** if  $\text{Im } C$  is a linear subspace of  $\Sigma_{out}^n$ .

Before we proceed, we mention that we will adopt a common abuse of notation and consider  $C$  alternately as both a function and as the set of encoded codewords. That is to say, given  $m \in \{0, 1\}^k$ , we denote by  $C(m)$  an encoded codeword, and given  $x \in \{0, 1\}^n$ , we denote by  $x \in C$  the fact that  $x$  is a codeword (i.e. - there exists an  $m \in \{0, 1\}^k$  such that  $C(m) = x$ ).

Clearly, our goal is to maximize both  $\delta, \rho$ . A natural question to ask is: how good can our code be?

### 1.1.2 Lower and Upper Bounds

**Claim 3.** (*Singleton Bound*) For any code  $C$  with relative distance and rate  $\delta, \rho$  it holds that  $\delta + \rho \leq 1 + \frac{1}{n}$

*Proof.* Let  $C : \{0, 1\}^k \rightarrow \{0, 1\}^n$  be a code with rate  $\rho = \frac{k}{n}$  and distance  $\Delta = \delta \cdot n$ . Consider the code  $C'$  which is given by removing the first  $\Delta - 1$  coordinates of any codeword  $C(m)$ . As  $d(C(m_1), C(m_2)) \geq \Delta$ , the function  $\text{Trim} : C \rightarrow C'$  is injective, as if  $\text{Trim}(x) = \text{Trim}(y)$  for some  $x, y \in C$  then  $d(x, y) \leq \Delta - 1 < \Delta$ . In particular,  $|C'| = |C| = 2^k$ . On the other hand, by our construction  $C' \subseteq \{0, 1\}^{n-d+1}$ , thus clearly  $|C'| \leq 2^{n-d+1}$ . Together, we get:

$$2^k \leq 2^{n-d+1}$$

Taking log and dividing over  $n$  gives the claim

□

Thus we see that we have a natural limit on how good our parameters can be. However, unlike most cases in theoretical computer science where achieving some upper bound is a long sought after goal, there is a simple, explicit construction, that achieves this bound:

#### 1.1.2.1 The Reed-Solomon Error Correcting Code

**Claim 4.** For every  $\delta, \rho > 0$  s.t.  $\delta + \rho = 1 + \frac{1}{n}$ , there exists an explicit linear code with relative distance  $\delta$  and rate  $\rho$ .

*Proof.* (Reed Solomon Code) As our construction is an algebraic one, we will require an important algebraic tool, namely: the fundamental theorem of Algebra.

**Theorem 1.5.** (*The Fundamental Theorem of Algebra, FTA*) Let  $\mathbb{F}$  be a field and  $f \in \mathbb{F}[x]$  be some non-zero polynomial of degree  $d$  over  $\mathbb{F}$ , then  $f$  has at most  $d$  roots in  $\mathbb{F}$

With this, consider the following code  $C : \{0, 1\}^k \rightarrow \mathbb{R}^n$ : given a message  $m = m_0, m_1, \dots, m_{k-1} \in \{0, 1\}^k$ , we define the polynomial  $f_m(x) \stackrel{\text{def}}{=} \sum_{i=0}^{k-1} m_i \cdot x^i$  and we encode the message by evaluating  $f_m$  over  $n$  distinct points, i.e.

$$C(m) = (f_m(0), f_m(1), \dots, f_m(n-1))$$

We first note that it is easy to verify that this is indeed a linear codeword. This is a direct consequence of the fact that if  $f_{m_1}, f_{m_2}$  are two polynomials then

$$(\alpha \cdot f_{m_1} + \beta \cdot f_{m_2})(x) = \alpha \cdot f_{m_1}(x) + \beta \cdot f_{m_2}(x)$$

Now, given two distinct messages  $m_1 \neq m_2 \in \{0, 1\}^k$ , the distance between  $C(m_1), C(m_2)$  is the number of evaluation points  $i$  such that  $f_{m_1}(i) \neq f_{m_2}(i)$ . Equivalently, this is the number of non-zero evaluation points of the codeword  $C(m_1 - m_2)$ . We now invoke the FTA - as  $m_1 \neq m_2$ , we know that  $f_{m_1 - m_2} \not\equiv 0$  and that it is a polynomial of degree at most  $k-1$ , and thus has at most  $k-1$  roots in  $\mathbb{R}$ . This implies that  $f_{m_1 - m_2}$  has at least  $n - k + 1$  non zero evaluations over said range, and thus  $\delta \cdot n \geq n - k + 1$  which implies  $\delta + \rho \geq 1 + \frac{1}{n}$  as needed.

The only problem with the above construction is that the output alphabet of the code is big. Indeed, considering only the highest monomial in  $f_m, x^{k-1}$ , we have that  $f_m(n) \approx n^k$ , i.e., exponential in  $k$ .

To fix this issue, we transition our construction to the realm of finite fields. Recall that for any prime power  $q$  there exists a field  $\mathbb{F} = \mathbb{F}_q$  of  $q$  elements. We will pick  $n \leq q \leq 2n$  and further insist that  $q$  will be a prime number and not just a prime power (this is not a necessity but will ease the notation in our construction). We recall that such a prime exists by Bertrand's postulate, and that  $\mathbb{F} \cong \mathbb{Z}_q$ , i.e. our fields has the elements  $0, 1, \dots, q-1$  and addition and multiplication are defined modulu  $q$  (we will see later on in the course that this is indeed a field).

Given the above, we construct our new code  $RS : \mathbb{F}^k \rightarrow \mathbb{F}^n$  in the same exact way - given a message  $m_0, \dots, m_{k-1}$  we define  $f_m(x) = \sum_{i=0}^{k-1} m_i x^i \in \mathbb{F}[x]$ . As  $\mathbb{F}$  is a field, any non-zero polynomial has at most  $k-1$  roots in the field and the distance bound remains intact while the output alphabet remains "small" (note that we still require  $|\mathbb{F}| \geq n$ , but this is much better than the exponential dependency we had before)  $\square$

### 1.1.2.2 Remarks

One might ask why go through the trouble of working over finite fields, and not simply interpret and evaluate the polynomial over  $\mathbb{Z}_n$  for some arbitrary  $n$ . This is necessary as a crucial component of our construction was the use of the FTA. Indeed, if we

consider the polynomial  $f(x) = x^2 + x$  over  $\mathbb{Z}_6$ , one can easily verify that 0, 2 and 3 are roots of  $f$  while  $\deg(f) = 2$ .

However, requiring that the field we use has enough elements to evaluate a polynomial over  $n$  points, i.e.  $p \geq n$  is costly in itself. The code we get would be some function  $\text{RS} : \mathbb{Z}_p^k \rightarrow \mathbb{Z}_p^n$ . Note that as  $n$  grows to  $\infty$ , so does the size of our field (a disadvantage).

Finally, we performed the analysis of the minimal distance of the code not by assessing the actual difference between two codewords, but by comparing the distance of any arbitrary codeword to the zero vector, this is not a coincidence and is a useful feature of linear codes:

**Lemma 1.6.** *Let  $C$  be a linear codeword of minimum distance  $d$ , then the following holds:*

$$d = \min_{x_1 \neq x_2 \in C} d(x_1, x_2) = \min_{\bar{0} \neq x \in C} d(x, \bar{0})$$

*Proof.* As  $C$  is a linear code,  $\bar{0} \in C$  and thus  $\min_{\bar{0} \neq x \in C} d(x, \bar{0}) \geq d$ . On the other hand, if  $x_1, x_2$  are two distinct codewords that achieve the minimal distance then  $d(x_1, x_2) = d(x_1 - x_2, \bar{0}) = d$ . The lemma follows as the code is linear and thus  $x_1 - x_2 \in C$  □

## 1.2 Expander Graphs

We will present the problem of expander graphs in short. For different (and countless) use-cases, we want a graph that is:

1. Extremely well connected, i.e. - in order to isolate a large group of vertices we need to remove many edges from the graph
2. Sparse (without too many edges).

Consider the **clique-graph**  $K_n$ , which is obviously very well connected. If we want to remove a vertex from it, we need to erase  $n - 1$  edges from it. If we want to remove  $|S|$  vertices, we need to erase  $|S|(n - |S|)$  edges. This graph clearly achieves the connecteness requirement. However, it is as far away from sparse as possible, i.e. - for any vertex  $v$ ,  $\deg v = n - 1$

On the other side of the spectrum, consider an arbitrary **tree-graph**  $T_n$ . In order to remove a single vertex one need only remove a single edge and there are even cases where removing  $|S|$  vertices can be achieved by removing a single edge. On the other hand, this graph is clearly as sparse as possible (for a connected graph).

Finally, consider the **cycle-graph**  $C_n$ : this graph is still very sparse (it has  $n$  edges compared to the minimal  $n - 1$  edges needed to maintain a connected graph). This graph is slightly better than the tree-graph as one must remove 2 edges to separate the graph into two connected components. Better, but still very poor.

**Can we have a well connected graph that is not much bigger than the tree?**

Turns out we can. We will show a well connected graph with only  $3n$  edges.

**Definition 1.7.** Let  $G = (V, E)$  be an undirected  $d$ -regular graph (i.e. the degree of each vertex is  $d$ ). For two subsets  $S, T \subseteq V$  we define  $E(S, T) = \{(u, v) \in E : u \in S \wedge v \in T\}$ .

We say that  $G$  is **well connected** if for any  $S \subseteq V$  such that  $|S| \leq n/2$  we have  $|E(S, V - S)| \geq \alpha \cdot d \cdot |S|$  for some constant  $\alpha > 0$ .

Note that  $1 \cdot d \cdot |S|$  is the best one can hope for (in the case where  $E(S, S) = \phi$ ).

We now present a construction which gives an expansion factor of  $0.1 \cdot d \cdot s$ , which is great. While being very simple to describe, the proof of the construction relies on deep number theory tools. For more information, see for example Theorem 4.4.2 in ?.

### 1.2.1 The Construction

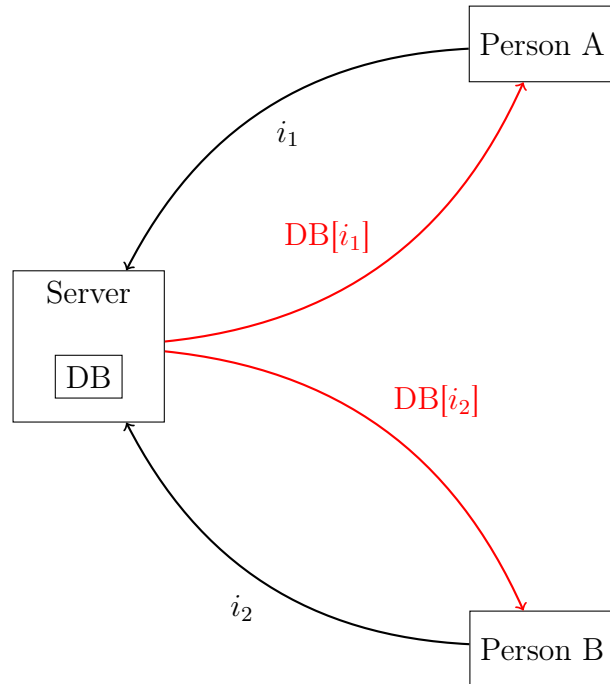
Define the graph  $G = (V, E)$  where we identify vertices with elements from  $\mathbb{Z}_p$ , and connect each node  $x \in \mathbb{Z}_p$  to 3 vertices:

1. Its multiplicative inverse  $x^{-1}$  (for  $x = 0$  we add a self-loop)
2. Its predecessor  $x - 1$
3. Its successor  $x + 1$

where all calculations are over  $\mathbb{Z}_p$ . Note that this graph is 3-regular, with  $3n$  edges. As stated, one can prove that this graph is an expander with  $0.1 \cdot 3 \cdot |S|$  crossing edges for every  $S \subseteq V$  of size at most  $p/2$  (but we will not show that). Note that  $p$  need not be a constant, and thus one can construct graphs of arbitrary size using this method.

## 1.3 Private Information Retrieval (PIR)

Consider the following scenario. We have a database and we want many people to access it. We can put it on one server - a naive construction.



**The Problem:** In the naive construction, the server knows  $i$ , the queried index, is connected to the user, and the user has no privacy. We're going to further require that the server learns **nothing** from the query.

**Definition 1.8.** A **communication protocol** between two parties  $A, B$  is a sequence of messages:

- $A$  sends a message  $a_1 \in \{0, 1\}^{m_1}$  to  $B$
- $B$  replies sending  $b_1 \in \{0, 1\}^{m_2}$  to  $A$
- $A$  sends  $a_2 \in \{0, 1\}^{m_3}$  to  $B$  and so forth

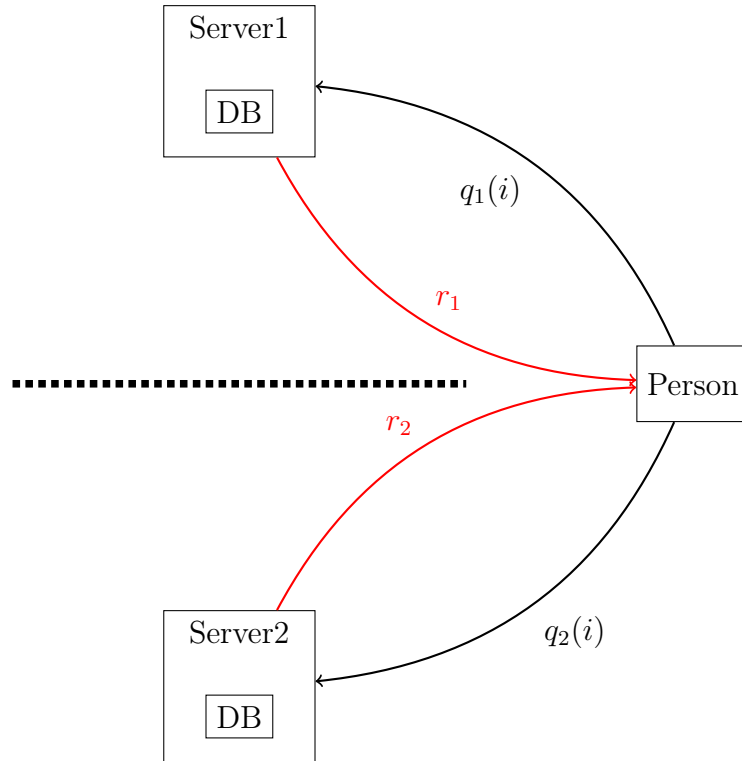
Given an execution of this protocol we call the sequence  $a_1, b_1, a_2, \dots$  the **transcript** of the protocol.

We say that a communication protocol between two parties  $A, B$  has **communication complexity**  $r$  if there exists a transcript between the two parties such that  $|a_1| + |b_1| + |a_2| + \dots = r$

A natural question is - what is the communication complexity under our requirements? One can show that in this model the only solution is for the server to send the the whole DB in response to a query, i.e., the query complexity is the size of the entire database.

### 1.3.1 Using Two Servers

It turns out we can do much better if we use two servers, that are not allowed to communicate.



Intuitively, think of server 1 returning  $b \oplus q_i$  and server 2 returning  $b$ , for some random independent bit  $b$  - obviously the two servers will learn nothing about  $q_i$ .

The original paper showed a construction for the two-server model that satisfies the privacy requirement, using  $r = O(n^{\frac{1}{3}})$  communication bits. Recently, a protocol was constructed which requires only  $r = O(2^{\sqrt{\log n}}) = n^{o(1)}$  bits of communication.

## 1.4 Randomness Mergers

### 1.4.1 The setup

Before we start, we give some basic definitions:

**Definition 1.9.** Let  $X, Y$  be a random variables over some universe  $\Omega$ , we define the following:

- The **support** of  $X$ :

$$\text{Sup}(X) = \{\omega \in \Omega : \Pr_{x \sim X}[x = \omega] > 0\}$$

- If  $X, Y$  have the same distribution we will denote this by  $X \sim Y$
- We define the **statistical distance** between the distributions:

$$\text{sd}(X, Y) = \frac{1}{2} \sum_{\omega \in \Omega} \left| \Pr_{x \sim X} [x = \omega] - \Pr_{y \sim Y} [y = \omega] \right|$$

We denote the case where  $\text{sd}(X, Y) \leq \epsilon$  by  $X \sim_{\epsilon} Y$  and we say that  $X$  is  $\epsilon$ -close to  $Y$

- Finally, in the case where  $\Omega = \{0, 1\}^k$ , we denote the uniform distribution over  $\Omega$  by  $U_k$

Consider the following scenario: we are given access to two random variables  $X, Y$  such that  $\text{Sup}(X), \text{Sup}(Y) \subseteq \{0, 1\}^n$  and:

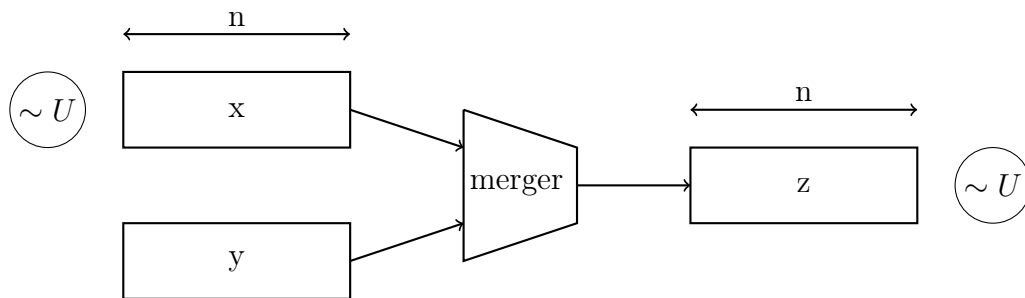
- Either  $X \sim U_n$  or  $Y \sim U_n$  (we don't know which)
- The random variables  $X, Y$  may be correlated

Our goal is to describe an algorithm, such that given  $X$  and  $Y$  as inputs, we produce a random variable  $Z$  such that, say,  $Z \sim U_{0.9n}$ , formally:

**Definition 1.10.** A function  $\text{Mer} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^{\ell}$  is called a **perfect randomness merger** with output length  $\ell$ , if given two (possibly correlated) random variables  $X, Y$  such that  $X \sim U_n$  or  $Y \sim U_n$ :

$$\text{Mer}(X, Y) \sim U_{\ell}$$

Where we define the output of  $\text{Mer}(X, Y)$  as follows: draw two instances  $x \sim X, y \sim Y$  and output  $\text{Mer}(x, y)$



Sadly, our next task is to show that no such function exists.



## 1.4.2 What cannot be achieved

**Claim 11.** *There are no perfect randomness mergers with output length 1*

*Proof.* Before we prove the claim, one might ask: why not xor the random variables, i.e.,  $\text{Mer}'(X, Y) = X \oplus Y$ ? It is not hard to check that this works in the case where  $X, Y$  are independent, but letting  $X \sim U_n$  and  $Y \equiv X$ , clearly  $\text{Mer}'(X, Y) \equiv 0$ . We now turn to prove the general case. Let  $F : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$  be an arbitrary function. Consider two cases:

- If there exists an  $x_0 \in \{0, 1\}^n$  such that  $x_0$  determines  $F$ , i.e.  $F(x_0, y) = b$  for any  $y \in \{0, 1\}^n$ , then we can let  $X \equiv x_0$  and  $Y \sim U_n$ . It is easy to see that  $F(X, Y) \equiv b$ , i.e.,  $F(X, Y)$  is constant
- If on the other hand for any  $x \in \{0, 1\}^n$  there exists  $y_0, y_1$  such that  $F(x, y_0) = 0, F(x, y_1) = 1$  then we do the following - define two functions  $Y_0, Y_1 : \{0, 1\}^n \rightarrow \{0, 1\}^n$  such that  $Y_0(x) = y_0, Y_1(x) = y_1$  and define  $X \sim U_n, Y = Y_0(X)$ . Again, it is easy to verify that  $F(X, Y) \equiv 0$

□

## 1.4.3 What can be achieved

Clearly, our hopes were too high. To achieve our goal we lower our expectations two-folds:

- First, we allow our merger some extra, independent randomness
- Secondly, we require that the output of the merger be  $\epsilon$ -close to a uniform distribution

Formally:

**Definition 1.12.** A function  $\text{Mer} : \{0, 1\}^n \times \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^\ell$  is called a **randomness merger** with output length  $\ell$  and error  $\epsilon$ , if given two (possibly correlated) random variables  $X, Y$  such that  $X \sim U_n$  or  $Y \sim U_n$ :

$$\text{Mer}(X, Y, U_d) \sim_\epsilon U_\ell$$

Where we define the output of  $\text{Mer}(X, Y, U_d)$  as follows: draw three instances  $x \sim X, y \sim Y$  and  $r \sim U_d$  (independently of  $X, Y$ ) and output  $\text{Mer}(x, y, r)$

Note that if  $d = n$  then one could trivially construct the randomness merger  $\text{Mer}(\cdot, \cdot, U_d) = U_d$ , thus our goal is obviously to construct such a function for  $d \ll n$ . In a fairly recent line of work it was shown that one can construct a randomness merger using  $d \approx 2 \log n$  extra bits of randomness. We now present the construction of such a merger. It is interesting to note that the correctness of the construction relies on the resolution of (the finite case of) a long standing conjecture in the field of geometric measure theory (see ?).

The randomness merger is constructed as follows:

1. Partition  $X$  and  $Y$  into  $\frac{n}{\log n}$  "blocks" of length  $\log n$  bits  $(X_1, \dots, X_{\frac{n}{\log n}}, Y_1, \dots, Y_{\frac{n}{\log n}})$
2. Choose a prime  $p$  s.t.  $\log p \approx \log n$
3. Choose random numbers  $A$  and  $B$ , independently and uniformly over  $Z_p - \{0\}$
4. Interpret the blocks  $X_i$  and  $Y_i$  as elements of  $Z_p$  and define  $Z_i = (A \cdot X_i + B \cdot Y_i) \bmod p$
5. Interpret each output  $Z_i$  as a binary string in  $\{0, 1\}^{\log n}$
6. Finally, output  $(Z_1, \dots, Z_{\frac{n}{\log n}})$

Interestingly, this simple and elegant construction is a randomness merger with error  $\epsilon = O\left(\frac{1}{n}\right)$ , as required.

## LECTURE 2

# THE FUNDAMENTAL THEOREM OF ALGEBRA; GROUPS - THE BASICS

---

In this lecture, we will discuss the Fundamental Theorem of Algebra while exploring the number systems  $\mathbb{N} \subseteq \mathbb{Z} \subseteq \mathbb{Q} \subseteq \mathbb{R} \subseteq \mathbb{C}$ . In retrospective, one can think of this process emerging so that the theorem will hold. The extension of some of these number systems from their prior will be abstracted later in the course and so it is beneficial to see the ideas involved on a familiar ground. I will mostly follow Chapters 1,2 of Stewart's excellent book on Galois Theory [Stewart \[2015\]](#).

## 2.1 The Familiar Number Systems

Solving polynomial equations, despite its boring reputation, has a fascinating history and required some significant psychological leaps from the very best of mathematicians (and it still does from the best of students). Slowly but surely, mathematicians extended their “number systems” when encountered with a problem expressed within the known number system whose solution was “outside” of it. In this section we briefly review this process. Some of the ideas that are required for extending these number systems will be abstracted later in the course and so it is beneficial to recall these ideas when applied at a familiar ground which we, at the very least, think we understand.

It all started with the set  $\{1, 2, 3, \dots\}$  which by itself is a completely non-trivial concept. It was highly abstract a few thousand years ago. It also didn't help that this set is infinite. Even today, many high school students are confused about the alleged paradox that every number is “finite” yet there are infinitely many of them. In this number system we can solve equations like  $x + 1 = 2$ . I will leave this as an exercise.

### 2.1.1 Enter zero

The acceptance of zero as a legitimate number took some getting used to. The ancient Greeks, for example, had no symbol for zero as they baffled with deep philosophical questions such as “how can nothing be something?”. Zero was used as a placeholder quite early in positional number systems like we use today, but it was considered nothing more for years to come. We write  $\mathbb{N} = \{0, 1, 2, \dots\}$  for the *natural numbers*.

Yes, we consider 0 to be a (very) natural number. In fact, when we come to formalize the notion of a number system using an axiomatic approach, the existence of (the abstraction of) 0 will be one of the axioms. More so, it will be the *only* number we demand to exist within the number system.

### 2.1.2 Negative numbers

Don't get me started about the negatives which allows one to solve equations like  $x + 1 = 0$ . It suffices to say that even at 1759, the English mathematician Maseres wrote that negative numbers "darken the very whole doctorines of equations and make dark the things which are in their nature excessively obvious and simple". Leibniz is considered to be the first to systematically employ negative numbers. He did so for his development of Calculus. I don't know about you, but I always imagined that Calculus is light years away from any discussion about negative numbers. Anyhow, denote the *whole numbers* by  $\mathbb{Z} = \{0, \pm 1, \pm 2, \dots\}$ . The letter  $\mathbb{Z}$  comes from the German word "Zahl" which translates to a teller in English.

### 2.1.3 Be rational

What about  $2x = 1$ ? Positive fractions seem to have been recognized earlier than zero and the negatives. However, there is some complexity involved in their formal definition. We are used to think of rational numbers as, well, numbers or more precisely as a pair of whole numbers. In particular, we write  $\mathbb{Q} = \{\frac{a}{b} \mid a, b \in \mathbb{Z}, b \neq 0\}$ . However, we identify some of the numbers in this set such as  $\frac{1}{2}$  and  $\frac{2}{4}$ . So, in fact, a rational number is not quite a pair of  $\mathbb{Z}$  elements but rather a set of such pairs. More precisely, a rational number is an equivalent class with respect to some equivalent relation. However, we are so used to this that we suppress this fact and, in particular, write things like  $\mathbb{Z} \subset \mathbb{Q}$  which formally does not make much sense. What we actually mean is that there is a copy of  $\mathbb{Z}$  "embedded" in  $\mathbb{Q}$ . This copy is given by  $\{\frac{a}{1} \mid a \in \mathbb{Z}\}$  and it behaves like  $\mathbb{Z}$  when we add and multiply unlike, say,  $\{\frac{1}{a} \mid a \in \mathbb{Z}\} \cup \{0\}$ .

*Exercise 1.* The Egyptians only considered fractions of the form  $\frac{1}{a}$  for  $a \in \{1, 2, 3, \dots\}$  (and  $\frac{2}{3}$  but let's ignore that one). One nice and not completely trivial fact is that any fraction  $\frac{a}{b}$  with  $1 \leq a \leq b$  can be written as a finite sum of distinct Egyptian fractions. Can you prove that?

Later in the course we will abstract this process of taking a number system like  $\mathbb{Z}$ , some of whose elements cannot be inverted, and "embed" it in a bigger number system that is closed to inversion.

### 2.1.4 Get real

Attempts made by ancient mathematicians who recognized only  $\mathbb{Q}$  as the set of numbers to solve  $x^2 = 2$  is a famous story in the history of Mathematics. Once again, the realization that this is impossible came as a philosophical shock.

*Exercise 2.* Here is a lesser-known proof sketch for the insolvability of  $x^2 = 2$  in  $\mathbb{Q}$ . Try to fill in the details. Assume by way of contradiction that  $\frac{a}{b}$  is a solution to  $x^2 = 2$  with  $a, b \in \mathbb{N}$  and  $b$  minimal among all such solutions. Consider now the expression  $\frac{2b-a}{a-b}$ .

Extending  $\mathbb{Q}$  to  $\mathbb{R}$  is completely non-trivial. It involves taking the topological closure of  $\mathbb{Q}$  with respect to the natural metric and by that close the (many many) “holes” in  $\mathbb{Q}$ . This is more or less done by adjoining the limits of all convergent sequences in  $\mathbb{Q}$ . Anyhow, whatever  $\mathbb{R}$  is, it is fairly safe to say that we all feel comfortable with it. We don’t call them real numbers for nothing!

### 2.1.5 Complex numbers

What about solving  $x^2 + 1 = 0$ ? We are all programmed to shout  $i$  (or  $\pm i$ ) but deep inside one might have the feeling that  $i$  is just a made up symbol—a cheat if you will. I mean,  $\sqrt{2}$ , I can get—it is the limit of a sequence of approximate solutions to  $x^2 = 2$ . But  $i$  is just, well, not real... Like the zero and the negative numbers,  $i$  wasn’t greeted with a smile by humankind. It was more like, well, we really need this guy to solve equations, but it was considered as this formal symbol that one can manipulate but dare not consider as “real”.

Let’s elaborate on that. We all know how to solve the general quadratic equation  $ax^2 + bx + c = 0$ . We have this neatly wrap expression for the solutions

$$x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

which I’m positive you can cook up by yourself. This formula, expressed quite differently, was known already to the Babylonians some 3600 years ago. Applying this to  $x^2 + 1 = 0$  doesn’t give any meaningful answer in  $\mathbb{R}$  as the  $\sqrt{\cdot}$  is applied to a negative number. This wasn’t a problem to  $i$ -non-believers. For them, it was simply Math’s way of telling us that there is no solution.

$i$  came to hunt the human race when people were finally able to solve cubic equations. It turns out that there is a general solution to such equations and one can derive it in a page or two (see Stewart’s book). However, there is a significant amount of trickery involved and it was an open problem to come up with a solution for quite some time. It was only at around 1535 that the general cubic equation was solved by

Fontana (nicknamed Tartaglia). First, using some standard trickery, one can reduce the general cubic equation to the form  $x^3 + px + q = 0$ . A general solution is then given by, get ready for this,

$$x = \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} + \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}.$$

Impressive no doubt. But, here is the catch. If we apply this to  $x^3 - 15x - 4 = 0$  which clearly has a solution  $x = 4$ , we get  $x = \sqrt[3]{2 + \sqrt{-121}} + \sqrt[3]{2 - \sqrt{-121}}$ . Where is our beloved 4? Turns out that if you are willing to consider  $i$  as if it was a legitimate number, assuming all rules of arithmetics apply to him, you can extract 4 out of this mess.

You see, it is not just that  $x^3 - 15x - 4$  has solutions outside of the reals which you may or may not choose to consider as real. It is that even if these solutions are real as 4, our way of finding them gets out of  $\mathbb{R}$  before landing back safely. You might not be so impressed. After all, this is just one way of finding a solution. Perhaps the undesired visit of  $i$  is due to the algorithm (the formula) not the problem itself. Well, turns out that one can prove, in some formal sense, that any solution that is expressed by radicals (square roots, cubic roots, etc) will go through  $i$  even in some cases in which all roots are real. Indeed, Mathematics is trying to tell us something... Soon enough we'll start talking about "field extensions". The uncomfortable feeling we may have had with  $i$ -adding this artificial solution-will come to hunt us again. So, we better surface these feelings at a familiar ground.

At any rate, we define  $\mathbb{C} = \{a + ib \mid a, b \in \mathbb{R}\}$  where addition and multiplication are given by "extending" these operations from  $\mathbb{R}$  together with the rule  $i^2 = -1$ . So, multiplication is given by

$$\begin{aligned} (a + bi)(c + di) &= ac + adi + bci + bdi^2 \\ &= (ac - bd) + (ad + bc)i. \end{aligned}$$

Going back to our friend,  $\sqrt{2}$ . Come to think of it, if something is not real then it is  $\sqrt{2}$ . I mean it is an endless pattern-less string of digits. There is not enough atoms in the universe to represent this idealized number. So, I claim you have never seen the *real*  $\sqrt{2}$  in your life!

## 2.2 Will this ever end?

One of the many cool features of  $\mathbb{C}$  is that it is the end of this game.  $\mathbb{C}$  has the remarkable property that *any* polynomial equation with coefficients in  $\mathbb{C}$  has *all* of its solutions in  $\mathbb{C}$ . That's a great deal! We added only this single weird symbol  $i$  so as to obtain/invent/discover, you choose, a solution to the specific simple equation  $x^2 + 1 = 0$  and what I'm saying is that by doing that, we got all solutions to all polynomial equations for free even if the coefficients have  $i$ 's in them! Later in the course we will refer to number systems that have this property *algebraically closed*. When I say that  $\mathbb{C}$  is the end of the game, I don't mean that  $\mathbb{C}$  is the only number system with this property. I mean that it is the only one if you start from  $\mathbb{R}$ .

This property of  $\mathbb{C}$  is given by The Fundamental Theorem of Algebra. To state it, recall that a solution to a polynomial equation  $p(x) = 0$  is called a *root* of  $p$ .

**Theorem 2.3** (The Fundamental Theorem of Algebra). *A non-constant polynomial with coefficients in  $\mathbb{C}$  has a root in  $\mathbb{C}$ .*

From **Theorem 2.3** one can deduce that a degree  $n \geq 1$  complex polynomial has exactly  $n$  roots. Some of these roots though may repeat more than once. For example,  $x^2 - 2x + 1$  can be written as  $(x - 1)^2$  from which one would agree that 1 counts as 2 roots of the polynomial, whatever that means.

**Theorem 2.3** wasn't obvious even for the great mathematicians of the time. For example, Bernoulli proposed a counterexample of degree 4. The great Euler proved him wrong in a letter to Goldbach. Euler claimed he has a proof for all degrees  $n \leq 6$ . A proof for the general case had to wait for Gauss who used trigonometric series in his 1799' proof.

For the Ph.D. students who are reading this, Gauss proved the theorem while being a Ph.D. student. Just saying :) Gauss, being Gauss, subsequently gave 3 other proofs. By now, there are many proofs, none of which is very easy, but you can fit one to a page or two (see Stewart's book). I'm not going to give a proof here. I'm gonna do something even better—I'm gonna show you *why* the theorem is true! The proof sketch is "topological" in nature, i.e., we're going to stretch continuous stuff in a continuous way. Also, I am kind of going to assume that you know about polar presentation.

*Proof Sketch.* Say \* you are looking at a polynomial  $p(x) = a_0 + a_1x + \cdots + a_nx^n$  with  $a_n \neq 0$ . If  $a_0 = 0$  then  $x = 0$  is clearly a root of  $p$ . So,  $a_0 \in \mathbb{C}$  sits somewhere in the complex plane away from the origin. Consider the following thought experiment. Fix a real number  $r \geq 0$  and consider the circle of all  $x \in \mathbb{C}$  with modulus  $|x| = r$ . Where does  $p$  map this circle to? Well, I don't quite know. But, if  $r$  is very large (compared

---

\*When turned into a formal proof, replace with "Let  $p(x)$  be..."

Figure 1: A traversal in  $\mathbb{C}$

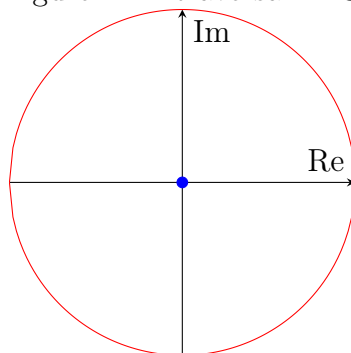
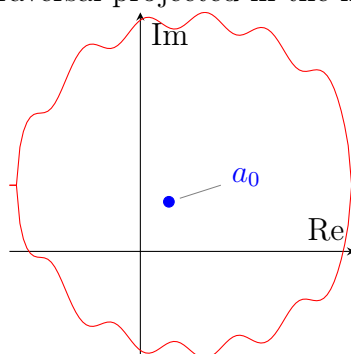


Figure 2: Traversal projected in the image of  $p(x)$



to the coefficients and the degree  $n$ ) then  $a_n x^n$  will be the dominating term. If it was the only term then the image of the map of the circle would have been a circle (in fact,  $n$  circles on top of each other) with modulus  $a_n r^n$ . However, there are these other pesky terms which make the actual image look like a wiggly circle. At any rate, if  $r$  is large enough the image is almost a circle.

Now comes the punch line. Starting from the huge  $r$  you chose, start to decrease it slowly all the way down to 0. If  $r$  is chosen large enough, we can make sure that the wiggly shape will contain both  $a_0$  and the origin. However, we know that at the end, when  $r = 0$ , the wiggly shape will converge to the single point  $a_0$  and, in particular. As everything we do is “continuous” at some point the wiggly shape—the image of  $p$ —must pass through the origin.  $\square$

## 2.3 What else is cool about $\mathbb{C}$ ?

Well, many things. For one, it turns out that  $\mathbb{C}$  is very real. I am no expert, but it seems that complex numbers are at the very least most suitable for describing Quan-



tum Mechanics. Mathematicians like complex numbers partially because working with complex functions is much nicer than with real-valued functions. To give some feeling for it, if you're working with a real function and it has an annoying singularity at some point, in  $\mathbb{R}$  the function is "broken" into two pieces. Over  $\mathbb{C}$  however you can just "go around" the misbehaved point. You can do much more though. For example, you can take a function that is defined somewhere in the complex plane but not in other possibly huge parts of it and, if the function is nice enough, you can extend it to more or less the whole complex plane in a unique way. It is a typical scenario that the new function shed new light on the original, partially defined, function. One fascinating application of such technique is to number theory and in particular to the Riemann Zeta Function. We'll talk a bit about it later in the course.

## 2.4 What does a Turing Machine think of $\mathbb{C}$ ?

The Fundamental Theorem of Algebra is extremely useful in theoretical computer science, coding theory, cryptography and what have you. However, computers (or Turing Machines if you must) don't like these infinite precision kind of number systems like  $\mathbb{R}$  and  $\mathbb{C}$ . Even  $\mathbb{Q}$  and  $\mathbb{Z}$  are not comfortable computing over as when turning to the analysis, one would need to keep track of the size of the computed numbers which, at best, is daunting.

Luckily, there are "finite number systems" which, being finite, avoid these issues. One can compute over these finite number systems and prove theorems about them. In particular, The Fundamental Theorem of Algebra more or less holds for these number systems as well—not just over  $\mathbb{C}$ . The proof, however, as you might expect looks very different as we're working in a very different setting. Soon we will get to these mysterious finite number systems. We will call them *finite fields*.

## 2.5 Bézout's Theorem

Another very interesting and useful generalization of [Theorem 2.3](#) is obtained by viewing the whole thing geometrically. First, let's work only over  $\mathbb{R}$  so it will be easier to draw things in our head. [Theorem 2.3](#) implies that over  $\mathbb{R}$ , a degree  $n \geq 1$  polynomial  $p(x)$  has at most  $n$  roots. Geometrically, this means that the set of points  $C = \{(x, p(x)) \mid x \in \mathbb{R}\}$  that describe the graph of  $p(x)$  in the real plane intersects the  $x$ -axis  $\{(x, 0) \mid x \in \mathbb{R}\}$  in at most  $n$  points. You can easily convince yourself that this holds true not only for the  $x$ -axis but actually for any line  $\{(x, y) \mid ax + by = c\}$  where  $a, b, c \in \mathbb{R}$ , not all zero, as long as  $C$  does not fully contain the line (this reservation with respect to the  $x$ -axis is hidden in the hypothesis of [Theorem 2.3](#) that

the degree  $n$  of  $p$  is greater than 1. This takes out the zero polynomial, whose graph is the  $x$ -axis, out of the picture.

We call  $C$  an *algebraic curve* (or simply a *curve*). Naturally, we say that  $C$  has degree  $n$ . The curves that correspond to linear equations have degree 1. So, **Theorem 2.3** implies that the number of intersection points between a degree  $n$  curve and a degree 1 curve in the plane is at most  $n \cdot 1$ . What about a degree  $n$  curve and a degree  $m$  curve? You guessed right! The number of intersection points is at most  $n \cdot m$ . This holds for even more general curves than “just” those of the form  $y = p(x)$ . You can mix up  $x, y$  in anyway you like. For example,  $xy - 1 = 0$  is a degree 2 curve.

This remarkable result is called *Bézout’s Theorem*. In fact, more is true. If you work over  $\mathbb{C}$  and count repeated points of intersections correctly you can almost say that the number of intersection points will be exactly  $n \cdot m$ . That is not quite true—think of two parallel lines. Turns out, though, that if you are open about changing your geometry from the standard geometry (called affine geometry) to what is called projective geometry, you get precisely  $n \cdot m$  points of intersection. The projective plane can be thought of as adding “points at infinity” to the affine plane, in which parallel lines meet. We won’t get into this in this course (ad ahead!) but will do so in a followup course on the fascinating subject of Algebraic Geometric codes. One of the goals of this course is to prepare you for the next one.

Let’s close with a fun fact. In his original 1770 paper, Bézout didn’t correctly account for multiplicities. As the theorem statement was “in the air”, one may argue (as some critics have) that the result is neither original nor correct...

## 2.6 Quick Introduction to Group Theory

In the previous lecture we used the informal notion of a “number system” while having  $\mathbb{N}, \mathbb{Z}, \mathbb{Q}, \mathbb{R}$  and  $\mathbb{C}$  in mind. I also hinted that we’ll be interested in *finite* number systems which clearly none of the above are. In this lecture we’ll start the beautiful process of formalizing and abstracting the notion of a number system. An abstraction that, in particular, will allow us to come up with finite number systems.

So what do we mean by a “number system”? Well, we definitely gonna want to have a set of “numbers” to play with. In all of the above examples, “playing” meant that we can add, subtract, multiply and, in some cases, even divide two numbers to obtain a third number in the set. A “number” in  $\mathbb{N}$  for example was a natural number. We have gotten used to those but to rigourously define such numbers one typically resorts to Peano axioms. For example, a set-theoretic model of the natural numbers, proposed by John von Neumann, constructs the natural numbers using only set operations on the empty set. You think you know good old 3?! According

to this construction, 3 is defined by  $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}\}$ , where  $\emptyset$  is the empty set. The situation got more complicated as we moved forward. In particular, a rational number is in fact an equivalence class with respect to a certain equivalence relationship over pairs of “numbers” from  $\mathbb{Z}$ .

When coming to abstract the notion of a number system, we don’t want to get into the details of what the numbers are. In particular, we won’t get into questions like what do they mean and are they “real” or not. We want to consider the numbers as abstract symbols and focus on how we can playing with them.

To make life simpler, in this lecture we are going to focus on a single operation rather than on four as we have in some of our number systems. This will make life simpler for us but in fact there are many examples, some of which we will see, in which the number system is naturally equipped with only one operation.

## 2.7 The definition of a Group

Defining our single-operation number system will be done using an axiomatic approach. Looking back at the number systems that we know, we’ll decide which properties we want to keep. Turns out we will set our heart on  $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}\}$  such properties.

### 2.7.1 The associative law

Let us consider  $\mathbb{Z}$  endowed with the addition operation, ignoring multiplication for the moment. One feature that we really like about addition is that although it is defined over pairs of numbers, we can extend the operation to more than two operands without having to account for the order in which we apply the operation. That long sentence can be summarized by saying that for all  $a, b, c \in \mathbb{Z}$ ,  $(a + b) + c = a + (b + c)$ . This “law” is known as the *associative law*. It implies that the expression  $a + b + c$  is well-defined. So does  $a + b + c + d$ , etc. We’re definitely going to keep this one! Note that not all operations are associative. We don’t have to look far:  $(a - b) - c$  is not the same thing as  $a - (b - c)$ .

### 2.7.2 Neutral element

OK. Now it is time to fix the poor treatment zero got from our ancestors and consider it as central. 0 “does nothing” when added to another element. We are going to insist our number systems will always have such a does-nothing element. Note that with respect to multiplication, 1 is the does-nothing element. Soon we are going to call such an element *a neutral element* rather than a does-nothing element.

### 2.7.3 Inverses

It is the negatives' turn. Still with  $\mathbb{Z}$  equipped with addition, every element  $a \in \mathbb{Z}$  has an element that when added to  $a$  gives back 0—the neutral element. Of course, we are talking about  $-a$ . When thinking of  $\mathbb{Q}$  equipped with multiplication the corresponding element to  $a$  will be  $1/a$ . But something about that example doesn't sit well—there is no “inverse” element like that which corresponds to 0. We will insist on having such an inverse element for *every* element. What is typically done in situations as in the above example is to shamefully take 0 out of the picture and consider  $\mathbb{Q} \setminus \{0\}$ . Sorry 0. Sometimes though it is useful and more “complete” to add  $\infty$  as the inverse of 0 (this is related to the projective geometry I eluded to last time). Geometrically, in the  $\mathbb{Z}$  with addition example, inverse means reflecting around 0 on the real axis. Taking, the inverse in  $\mathbb{C} \setminus \{0\}$  with respect to multiplication kind of turns the unit ball  $\{a \in \mathbb{C} \setminus \{0\} \mid |a| \leq 1\}$  inside out (and in some sense, if insisting, taking the omitted point 0 very far to  $\infty$ ).

### 2.7.4 The formal definition

After the motivating discussion above, I think we're ready to see the following definition.

**Definition 2.4** (Group). A *group* is a set  $G$  together with a function  $f: G \times G \rightarrow G$  such that

**Associativity.**  $\forall a, b, c \in G \quad f(f(a, b), c) = f(a, f(b, c));$

**Neutral element.**  $\exists e \in G \quad \forall a \in G \quad f(a, e) = f(e, a) = a;$

**Inverse.**  $\forall a \in G \quad \exists b \in G \quad f(a, b) = e.$

The function  $f$  is called the *group law*. We typically do not use this functional notation and simply write  $a \bullet b$  for  $f(a, b)$ . So, associativity looks like  $(a \bullet b) \bullet c = a \bullet (b \bullet c)$ . When it causes no confusion, one omits the symbol altogether and simply write  $ab$ . When we have an additive operation in mind we write  $a + b$  even though this is not necessarily the addition we know from  $\mathbb{N}$  and friends. In this case we sometime write 0 for  $e$ . We say that the group is *additive*. Similarly, when we have a multiplicative operation in mind, we write  $a \cdot b$  and 1 for  $e$ . In this case we say that the group is *multiplicative*. In an additive group one abbreviates and write  $a + a$  as  $2a$ ,  $a + a + a$  as  $3a$  etc. Similarly, in a multiplicative group we write  $a^2$  for  $a \cdot a$ , etc. Note that then it makes sense to define  $a^0 = e$ . It is important to keep in mind that these are just notations. They mean nothing and are only meant so that us human beings will convey our intensions (somewhat like comments in a code).

When referring to a group, we sometimes write  $(G, f)$ ,  $(G, \bullet)$  or  $(G, +)$ ,  $(G, \cdot)$ . When the operation is clear from context we sometimes write  $G$  for the group.

The innocent-looking [Definition 3.1](#) is central to a fair fraction of mathematics. It is hard to overestimate its importance. I will skip the historical development and nice anecdotes this time but I urge you to look it up. I will be satisfied by saying that the axioms were first written down formally by Walther von Dyck in 1882. However, the study of group theory, even before the definition has emerged, dates at least a century back.

### 2.7.5 But why this definition?

Chess is an interesting game. My eight year old son beats me more than half the time and I am not sure it says much about his game skills. Still, only by playing Chess I came to appreciate the rules—the axioms if you will—of the game. Prior to this experience, it all seemed a bit ad hoc.

You just saw [Definition 3.1](#)—the rules of the game. It is too much, to say the least, to expect of you to appreciate the definition just by staring hard at it. For that you will need to play with groups. We will only need the very basics of group theory but even that, I hope, will give you some insights. I will briefly mention that there are other notions related to groups. For example, it is worth noting that  $(\mathbb{N}, +)$  is not a group as none of the elements, but for 0, has an inverse. Such a structure is called a *monoid*—a group without the existence of inverse axiom. A *semigroup* does not even require the neutral element—only associativity. We won't touch upon these, somewhat less central, notions. Somehow, the structure guaranteed by the properties of a group is just enough to be extremely interesting. With hindsight, removing any of the axioms damages the structure and deem the resulting object less interesting. In the next section, we will go in the other direction and will add a “bonus” axiom to a group.

## 2.8 Commutative groups

The groups  $(\mathbb{Z}, +)$ ,  $(\mathbb{Z}, \cdot)$ ,  $\dots$ ,  $(\mathbb{C}, \cdot)$  all share one extra property that we did not insist on in the definition of a group. Can you see what it is? Yep, it is that the order of the operands does not change the result.  $2 + 3 = 3 + 2$  and  $2 \cdot 3 = 3 \cdot 2$ . We will care both about groups with this property and about groups without this property.

**Definition 2.5** (Commutative groups). A group  $(G, \bullet)$  is *commutative* if  $\forall a, b \in G \quad a \bullet b = b \bullet a$ .

Sometimes we use the term *Abelian group* for a commutative group named after the tragic genius Niels Henrik Abel.

An important example of a non-commutative group is obtained by considering composition of functions. Here, the set  $G$  consists of, say, all one-to-one functions on some domain  $D$ . The group law is then given by composition, namely,  $a \bullet b = a \circ b$  where, as usual,  $a \circ b$  is the function that is defined as follows:  $\forall x \in D \quad (a \circ b)(x) = a(b(x))$ . Convince yourself that this is a group and show that if  $|D| \geq 3$ , this group is non-commutative.

## 2.9 Constructing some groups

Let's try to play a little bit with groups of a given size to see what kind of structure the axioms dictate. As usual, it is good to start at the beginning. Is there a group of size 0? Well, nope. We insisted on having the neutral element and the empty set contains no elements. It is easy to see that there is exactly one group of size 1.

What about groups of size 2? Say  $G = \{e, a\}$  where  $e$ , as usual, is the neutral element. The axioms dictate how to multiply by  $e$  ( $ee = e$ ,  $ea = ae = a$ ). So, we only need to explore the two options  $aa = e$  and  $aa = a$ . If we go with the second option,  $a$  will not have an inverse as  $ae = aa = a$ . This contradicts the existence of an inverse axiom. You can convince yourself that the first option is consistent with all group axioms. So, there is exactly one group of size 2. In fact, as a computer scientist, you know that group very well. Do you recognize it? It is addition modulo 2—the way we are used to working with bits. This becomes more transparent if we choose to represent the group as an additive group, writing 0 for  $e$  and 1 for  $a$ . We denote this group by  $(\mathbb{Z}_2, +)$ .

One should be careful about what does it mean for two groups to be different. It is quite often the case that you're looking at this neat group you found just to realize it is a known group in disguise. That's one of the nice aspects of group theory. Just a moment ago, we figured out that our size 2 group is something we all know very well and might incorrectly consider it as being different. In the next lecture we will formalize what does it mean for two groups to be “the same”. But, for now, here is an example of two different, by all accounts, groups of size 4 (as an exercise, work out yourself all groups of size 3).

First, we have  $(\mathbb{Z}_4, +)$  - the group of addition modulo 4 over  $\{0, 1, 2, 3\}$ . Convince yourself this is indeed a group. Second, if we will “glue two independent” copies of  $(\mathbb{Z}_2, +)$  next to each other, I claim, we will be looking at a different group of size 4. Let me formalize what I mean by gluing. First, we recall the following notation - if  $A, B$  are two sets, we define  $A \times B = \{(a, b) \mid a \in A, b \in B\}$  as the product set.

**Definition 2.6** (Direct product). Let  $(G, \bullet_G)$ ,  $(H, \bullet_H)$  be two groups. We define the group  $(G \times H, \bullet)$  as follows. For every  $(g, h), (g', h')$  in  $G \times H$ , define  $(g, h) \bullet (g', h') = (g \bullet_G g', h \bullet_H h')$ .

The definition assumes implicitly that the resulting structure is a group. Convince yourself that this assertion indeed holds. One typically writes  $G \times H$  for the direct product between groups and omit the operation symbol. Note that its neutral element is  $(e_G, e_H)$  where  $e_G, e_H$  are the neutral elements of  $G, H$ , respectively. Note also that if  $G, H$  are finite, then  $|G \times H| = |G||H|$ .

With the direct product in hand, consider  $\mathbb{Z}_2 \times \mathbb{Z}_2$  which is a group of size 4. This is actually one famous group named the *Klein four-group*. There are many ways to convince yourself that this is not at all  $\mathbb{Z}_4$ . For example, every element in  $\mathbb{Z}_2 \times \mathbb{Z}_2$  is its own inverse whereas in  $\mathbb{Z}_4$ , 1 is not its own inverse (as  $1 + 1 = 2 \neq 0$ ). Turns out that there are no more groups of size 2 but to say such a thing formally we'll need to wait for the next lecture. Still, for now, here is a table of the number of groups of a given size up to 16. We also count commutative and non-commutative groups separately to get some more insight.

Size	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Commutative	1	1	1	2	1	1	1	3	2	1	1	2	1	1	1	5
Non-commutative	0	0	0	0	0	1	0	2	0	1	0	3	0	1	0	9

By the above examples, you have probably figured out by now that for every integer  $n \geq 1$  there is a group of size  $n$ . It is the group of addition modulo  $n$  over the set  $\{0, 1, \dots, n-1\}$ . We denote this group by  $(\mathbb{Z}_n, +)$ . For every  $n$ , this is a commutative group. Observe further that according to the table, for every prime number  $p \leq 16$ ,  $(\mathbb{Z}_p, +)$  is the only group of size  $p$ . This is true in general and we will prove this soon. Looking at size 15 one can see that this is not an if and only if condition

## 2.10 Some basic general properties

After seeing some groups we are ready to go abstract again and ask what kind of general properties hold for an abstract, general, group. First, we observe that a neutral element in every group is *unique*. Indeed, if  $e_1, e_2$  are both neutral then  $e_1 = e_1 e_2 = e_2$ . This allows us to talk about *the* neutral element in a give group.

An inverse of a given element is also unique. In fact, more is true. Consider a group  $G$  and a fixed element  $a \in G$ . Define the map  $\phi_a : G \rightarrow G$  by  $\phi_a(x) = ax$ . I claim that  $\phi_a$  is injective (that is, one to one). Indeed, if  $\phi_a(x) = \phi_a(y)$  then  $ax = ay$ . Let's take an inverse  $b$  of  $a$ . Then  $b(ax) = b(ay)$  and so, by the associativity law,

$(ba)x = (ba)y$ . As  $b$  is an inverse of  $a$ , we get that  $ex = ey$  and so  $x = y$ . To see that this implies the uniqueness of an inverse of a given element, consider some element  $a$ . As  $\phi_a(x)$  is injective, there is at most one element  $x$  for which  $ax = \phi_a(x) = e$  and, by the group axioms, we are guaranteed that such an element  $x$  exists. From now on we can talk about *the* inverse of an element  $a$ . If we have a multiplicative group in mind, we denote the inverse of  $a$  by  $a^{-1}$ . For an additive group we write  $-a$ . If we have a multiplicative group  $(G, \cdot)$  in mind,  $a \in G$  and  $n \in \mathbb{N}$ , you can prove that the inverse of  $a^n$  which is denoted by  $(a^n)^{-1}$ , equals to  $(a^{-1})^n$ . So we can actually extend the exponent to be an element of  $\mathbb{Z}$ . Same is true for an additive group.

## 2.11 Subgroups

In Mathematics, when you have an object that you care about that has some structure, you will typically be caring about parts of the object that also have such structure. Wow, that was abstract. Perhaps the simplest example is the integers  $(\mathbb{Z}, +)$  and even numbers as their sub-object. The structure of  $(\mathbb{Z}, +)$  is its group structure. The even numbers can also be added and they have a neutral element sitting inside. Associativity is obvious as it follows from the associativity of the integers. What we just said is that the even numbers, with respect to the addition operation is a group, which we'll denote by  $(2\mathbb{Z}, +)$ . But we really want to say that it is a *subgroup* of  $(\mathbb{Z}, +)$ . So, let's define this notion in the most natural way.

It will be useful to go back to the functional definition. We also use the following standard notation. Let  $A, B, C$  be three sets such that  $A \subseteq B$ . Let  $f: B \rightarrow C$ . We define the function  $f|_A: A \rightarrow C$  by  $f|_A(a) = f(a)$  for all  $a \in A$  and say that  $f|_A$  is the restriction of  $f$  to  $A$ .

**Definition 2.7** (Subgroup). Let  $(G, f)$  be a group. A *subgroup*  $(H, f|_{H \times H})$  of  $(G, f)$  is a group such that  $H \subseteq G$ .

Observe that as we require that  $(H, f|_{H \times H})$  is a group, we implicitly require that the image of  $f|_{H \times H}$  is contained in  $H$ .

A good example are vector spaces and their subspaces. For instance, take the group  $(\mathbb{R}^2, +)$ . Every line that intersects with the origin  $H = \{(x, y) | y = m \cdot x\}$  is a subgroup  $(H, +)$ . By the definition of the group  $(\mathbb{R}^2, +)$ , the subgroup holds all of the group's properties. We just need to show that  $\forall a, b \in H, a + b \in H$ . let  $(x_1, y_1), (x_2, y_2) \in H$ . so  $y_1 + y_2 = m(x_1 + x_2)$  so  $(x_1 + x_2, y_1 + y_2) \in H$ .



## 2.12 Cosets

We already saw that  $(2\mathbb{Z}, +)$  is a subgroup of  $(\mathbb{Z}, +)$ . Obviously  $(2\mathbb{Z} + 1, +)$  (i.e. the odd numbers) is not a subgroup, as it does not contain the neutral number 0. But still, it feels like it should be something like a subgroup, as the only action we did was to "move" the evens set by one step. We will call this type of set a *coset*. Notice that we can get the odds set by "moving" the evens set any odd number of steps, e.g.  $(2\mathbb{Z} + 3, +)$ . Actually, we will see that there are many ways to form the same coset. Let's start with the definition:

**Definition 2.8** (Coset). Let  $(G, f)$  be a group and  $(H, f)$  a subgroup ( $H < G$ ). Define  $aH$ , a *left coset* of  $H$ , to be the set  $\{ah | \forall h \in H\}$ . Similarly,  $Ha = \{ha | \forall h \in H\}$  is said to be a *right coset* of  $H$ .

*Example 2.9.* Considered again the group  $(\mathbb{R}^2, +)$ , a subset  $S = \{(x, y) | y = m \cdot x + n\} \subseteq \mathbb{R}^2$  isn't a subgroup of  $(\mathbb{R}^2, +)$  if  $n \neq 0$ , but it is a coset of  $H = \{(x, y) | y = m \cdot x\}$ , as  $S = (0, n) + H$ .

Some easy to prove claims:

1.  $G$  is commutative  $\implies \forall a \in G, \forall H < G \quad aH = Ha$
2.  $aH = bH \iff b^{-1}a \in H$
3.  $aH \neq bH \implies aH \cap bH = \emptyset$
4.  $|aH| = |H|$

Note that for a given subgroup  $H, \forall a, b \in G$  the relation  $aRb \iff b^{-1}a \in H$  defines an equivalence relation in  $G$ ; using (2), we can see that  $a$  and  $b$  are equivalent under this relation if and only if they belong to the same left coset of  $H$ .

We can conclude that the left cosets are a decomposition of the group  $G$  into disjoint sets of identical size:  $|G| = |\bigcup_a aH| = |H| \cdot \#\text{left cosets}$ .

**Definition 2.10** (index of  $H$  in  $G$ ). Let  $G$  be a group and  $H < G$  a subgroup. Then  $[G : H]$  the index of  $H$  in  $G$  is the number of distinct left cosets of  $H$  in  $G$ .

**Theorem 2.11** (Lagrange).  $|G| = |H| \cdot [G : H]$  i.e. the size of the group  $G$  can be expressed as the size of a subgroup  $H$  times the number of cosets of  $H$ . In particular, the size of the group  $G$  is divisible by the size of any subgroup of  $G$ .

Corollaries from Lagrange Theorem:

1. If  $G$  is a finite group,  $\forall H < G, |G|$  is divisible by  $|H|$ .
2. If  $G$  is prime then  $G$  doesn't have any non-trivial subgroups.

## 2.13 Normal groups and quotient groups

Consider a group  $G$  and a partition of  $G$  into equivalence classes under some equivalence relation  $R$ . We might want to define a group structure on the set of equivalence classes; that is, a group  $G'$  where each element is a set  $[a] = \{g \in G \mid aRg\} \subseteq G$ . We want  $G'$  to preserve the structure of  $G$ , so intuitively the new group operation should be defined in the following way:  $\forall [a], [b]$  elements in  $G'$ ,  $[a] \cdot [b] = [a \cdot b]$ . The problem with this approach is that this operation is not always well-defined: if  $[a] = [a']$ ,  $[b] = [b']$ , it is not necessarily true that  $[ab] = [a'b']$ .

*Remark.* Assume that  $G$  is a Group, and  $R$  is an equivalence relation such that it indeed holds that  $\forall a, a', b, b' \in G$  if  $[a] = [a']$  and  $[b] = [b']$ , then  $[ab] = [a'b']$ . In this case we could define a group  $G'$ , and:

1.  $[e]$  would be the neutral element in  $G'$ . proof:  $\forall [a] \in G'$   $[e] \cdot [a] = [e \cdot a] = [a] = [a \cdot e] = [a] \cdot [e]$
2.  $[e] = H$  is a subgroup of  $G$ :

**Associativity** :  $H \subseteq G$  therefore by the associativity of  $G$ ,  $H$  is also associative.

**Neutral element** :  $e \in H = [e]$ . Since  $e$  is the neutral element of  $G$ ,  $\forall h \in H \subseteq G$ ,  $e \cdot h = h \cdot e = h$ ;

**Inverse** :  $\forall h \in H$ ,  $\exists h^{-1} \in G$  by the inverse property of  $G$ .  $[h^{-1}] = [h^{-1} \cdot e] = [h^{-1}] \cdot [e] = [h^{-1}] \cdot [h] = [h^{-1} \cdot h] = [e]$ . This implies that  $h^{-1} \in [e]$ ;

**Closure under the group operation** :  $\forall a, b \in H$  :  $[ab] = [a][b] = [e][e] = [e] = H$ . This implies  $a \cdot b \in H$ .

3. Exercise: prove that  $[a] = aH = Ha$ .

We just saw that if we had some partition of  $G$  into disjoint sets that happen to have this nice property that  $\forall a, a', b, b' \in G$  if  $[a] = [a']$  and  $[b] = [b'] \implies [ab] = [a'b']$ , then  $[e]$  is a subgroup of  $G$ , the disjoint set of the partition are its left cosets, and they are identical to its right cosets. We can now understand the motivation for the following definitions:

**Definition 2.12** (Normal Subgroup). A subgroup  $N$  of a group  $G$  is normal in  $G$  if  $\forall g \in G : gN = Ng$ . We denote such a subgroup as  $N \triangleleft G$ .

Remarks:

1. In fact, the common definition for a normal subgroup is: a subgroup  $N$  of  $G$  is a normal subgroup of  $G$  if  $\forall g \in G$  and  $\forall n \in N, gng^{-1} \in N$ . It's easy to prove that the two definitions are equivalent.

2. Note that if  $G$  is abelian, then any subgroup of  $G$  is a normal subgroup.

**Definition 2.13** (Quotient Group). Let  $N$  be a normal subgroup of a group  $G$ . Define the quotient group  $G/N = \{aN : a \in G\}$ , i.e. the set of all cosets of  $N$  in  $G$ .  $G/N$  operation is defined as  $\forall aN, bN \in G/N : (aN) \cdot (bN) = (a \cdot b)N$

Note that  $G/N$  is indeed a group:

1. The operation is well defined as the definition of the product of two cosets does not depend on the representatives. If for some  $a, a', b, b' \in G : aN = a'N, bN = b'N$  then using  $N$  normality in  $G$  it holds that  $(ab)N = a(bN) = a(b'N) = a(Nb') = (aN)b' = (a'N)b' = a'(Nb') = a'(b'N) = (a'b')N$
2. The operation is associative:  $\forall aN, bN, cN \in G/N, (aNbN)cN = (abN)cN = ((ab)c)N = (a(bc))N = aN(bcN) = aN(bNcN)$
3.  $N = eN$  is the neutral element of  $G/N$  since  $\forall aN \in G/N, aN \cdot eN = (ae)N = aN = (ea)N = eN \cdot aN$
4.  $\forall aN \in G/N$ , there is an inverse in  $G/N$  and  $(aN)^{-1} = a^{-1}N: a^{-1}N \cdot aN = (a^{-1}a)N = eN = N$

## LECTURE 3

# GROUPS - THE BASICS

---

In the previous lecture we used the informal notion of a “number system” while having  $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{Q}$ ,  $\mathbb{R}$  and  $\mathbb{C}$  in mind. I also hinted that we’ll be interested in *finite* number systems which clearly non of the above are. In this lecture we’ll start the beautiful process of formalizing and abstracting the notion of a number system. An abstraction that, in particular, will allow us to come up with finite number systems.

So what do we mean by a “number system”? Well, we definitely gonna want to have a set of “numbers” to play with. In all of the above examples, “playing” meant that we can add, subtract, multiply and, in some cases, even divide two numbers to obtain a third number in the set. A “number” in  $\mathbb{N}$  for example was a natural number. We have gotten used to those but to rigourously define such numbers one typically resorts to Peano axioms. For example, a set-theoretic model of the natural numbers, proposed by John von Neumann, constructs the natural numbers using only set operations on the empty set. You think you know good old 3?! According to this construction, 3 is defined by  $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}$ , where  $\emptyset$  is the empty set. The situation got more complicated as we moved forward. In particular, a rational number is in fact an equivalence class with respect to a certain equivalence relationship over pairs of “numbers” from  $\mathbb{Z}$ .

When coming to abstract the notion of a number system, we don’t want to get into the details of what the numbers are. In particular, we won’t get into questions like what do they mean and are they “real” or not. We want to consider the numbers as abstract symbols and focus on how we can playing with them.

To make life simpler, in this lecture we are going to focus on a single operation rather than on four as we have in some of our number systems. This will make life simpler for us but in fact there are many examples, some of which we will see, in which the number system is naturally equipped with only one operation.

### 3.1 The definition of a Group

Defining our single-operation number system will be done using an axiomatic approach. Looking back at the number systems that we know, we’ll decide which properties we want to keep. Turns out we will set our heart on  $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}$  such properties.

### 3.1.1 The associative law

Let us consider  $\mathbb{Z}$  endowed with the addition operation, ignoring multiplication for the moment. One feature that we really like about addition is that although it is defined over pairs of numbers, we can extend the operation to more than two operands without having to account for the order in which we apply the operation. That long sentence can be summarized by saying that for all  $a, b, c \in \mathbb{Z}$ ,  $(a + b) + c = a + (b + c)$ . This “law” is known as the *associative law*. It implies that the expression  $a + b + c$  is well-defined. So does  $a + b + c + d$ , etc. We’re definitely going to keep this one! Note that not all operations are associative. We don’t have to look far:  $(a - b) - c$  is not the same thing as  $a - (b - c)$ .

### 3.1.2 Neutral element

OK. Now it is time to fix the poor treatment zero got from our ancestors and consider it as central. 0 “does nothing” when added to another element. We are going to insist our number systems will always have such a does-nothing element. Note that with respect to multiplication, 1 is the does-nothing element. Soon we are going to call such an element a *neutral element* rather than a does-nothing element.

### 3.1.3 Inverses

It is the negatives turn. Still with  $\mathbb{Z}$  equipped with addition, every element  $a \in \mathbb{Z}$  has an element that when added to  $a$  gives back 0—the neutral element. Of course, we are talking about  $-a$ . When thinking of  $\mathbb{Q}$  equipped with multiplication the corresponding element to  $a$  will be  $1/a$ . But something about that example doesn’t sit well—there is no “inverse” element like that which corresponds to 0. We will insist on having such an inverse element for *every* element. What is typically done in situations as in the above example is to shamefully take 0 out of the picture and consider  $\mathbb{Q} \setminus \{0\}$ . Sorry 0. Sometimes though it is useful and more “complete” to add  $\infty$  as the inverse of 0 (this is related to the projective geometry I was eluded to last time).

Geometrically, in the  $\mathbb{Z}$  with addition example, inverse means reflecting around 0 on the real axis. Taking, the inverse in  $\mathbb{C} \setminus \{0\}$  with respect to multiplication kind of turns the unit ball  $\{a \in \mathbb{C} \setminus \{0\} \mid |a| \leq 1\}$  inside out (and in some sense, if insisting, taking the omitted point 0 very far to  $\infty$ ).

### 3.1.4 The formal definition

After the motivating discussion above, I think we're ready to see the following definition.

**Definition 3.1** (Group). A *group* is a set  $G$  together with a function  $f: G \times G \rightarrow G$  such that

**Associativity.**  $\forall a, b, c \in G \quad f(f(a, b), c) = f(a, f(b, c));$

**Neutral element.**  $\exists e \in G \quad \forall a \in G \quad f(a, e) = f(e, a) = a;$

**Inverse.**  $\forall a \in G \quad \exists b \in G \quad f(a, b) = e.$

The function  $f$  is called the *group law*. We typically do not use this functional notation and simply write  $a \bullet b$  for  $f(a, b)$ . So, associativity looks like  $(a \bullet b) \bullet c = a \bullet (b \bullet c)$ . When it causes no confusion, one omits the symbol altogether and simply write  $ab$ . When we have an additive operation in mind we write  $a + b$  even though this is not necessarily the addition we know from  $\mathbb{N}$  and friends. In this case we sometime write 0 for  $e$ . We say that the group is *additive*. Similarly, when we have a multiplicative operation in mind, we write  $a \cdot b$  and 1 for  $e$ . In this case we say that the group is *multiplicative*. In an additive group one abbreviates and write  $a + a$  as  $2a$ ,  $a + a + a$  as  $3a$  etc. Similarly, in a multiplicative group we write  $a^2$  for  $a \cdot a$ , etc. Note that then it makes sense to define  $a^0 = e$ . It is important to keep in mind that these are just notations. They mean nothing and are only meant so that us human beings will convey our intentions (somewhat like comments in a code).

When referring to a group, we sometimes write  $(G, f)$ ,  $(G, \bullet)$  or  $(G, +)$ ,  $(G, \cdot)$ . When the operation is clear from context we sometimes write  $G$  for the group.

The innocent-looking [Definition 3.1](#) is central to a fair fraction of mathematics. It is hard to overestimate its importance. I will skip the historical development and nice anecdotes this time but I urge you to look it up. I will be satisfied by saying that the axioms were first written down formally by Walther von Dyck in 1882. However, the study of group theory, even before the definition has emerged, dates at least a century back.

### 3.1.5 But why this definition?

Chess is an interesting game. My eight year old son beats me more than half the time and I am not sure it says much about his game skills. Still, only by playing Chess I came to appreciate the rules—the axioms if you will—of the game. Prior to this experience, it all seemed a bit ad hoc.

You just saw **Definition 3.1**—the rules of the game. It is too much, to say the least, to expect of you to appreciate the definition just by staring hard at it. For that you will need to play with groups. We will only need the very basics of group theory but even that, I hope, will give you some insights. I will briefly mention that there are other notions related to groups. For example, it is worth noting that  $(\mathbb{N}, +)$  is not a group as none of the elements, but for 0, has an inverse. Such a structure is called a *monoid*—a group without the existence of inverse axiom. A *semigroup* does not even require the neutral element—only associativity. We won't touch upon these, somewhat less central, notions. Somehow, the structure guaranteed by the properties of a group is just enough to be extremely interesting. With hindsight, removing any of the axioms damages the structure and deems the resulting object less interesting. In the next section, we will go in the other direction and will add a “bonus” axiom to a group.

## 3.2 Commutative groups

The groups  $(\mathbb{Z}, +)$ ,  $(\mathbb{Z}, \cdot)$ ,  $\dots$ ,  $(\mathbb{C}, \cdot)$  all share one extra property that we did not insist on in the definition of a group. Can you see what it is? Yep, it is that the order of the operands does not change the result.  $2 + 3 = 3 + 2$  and  $2 \cdot 3 = 3 \cdot 2$ . We will care both about groups with this property and about groups without this property.

**Definition 3.2** (Commutative groups). A group  $(G, \bullet)$  is *commutative* if  $\forall a, b \in G \quad a \bullet b = b \bullet a$ .

Sometimes we use the term *Abelian group* for a commutative group named after the tragic genius Niels Henrik Abel.

An important example of a non-commutative group is obtained by considering composition of functions. Here, the set  $G$  consists of, say, all one-to-one functions on some domain  $D$ . The group law is then given by composition, namely,  $a \bullet b = a \circ b$  where, as usual,  $a \circ b$  is the function that is defined as follows:  $\forall x \in D \quad (a \circ b)(x) = a(b(x))$ . Convince yourself that this is a group and show that if  $|D| \geq 3$ , this group is non-commutative.

## 3.3 Constructing some groups

Let's try to play a little bit with groups of a given size to see what kind of structure the axioms dictate. As usual, it is good to start at the beginning. Is there a group of size 0? Well, nope. We insisted on having the neutral element and the empty set contains no elements. It is easy to see that there is exactly one group of size 1.

What about groups of size 2? Say  $G = \{e, a\}$  where  $e$ , as usual, is the neutral element. The axioms dictate how to multiply by  $e$  ( $ee = e$ ,  $ea = ae = a$ ). So, we only need to explore the two options  $aa = e$  and  $aa = a$ . If we go with the second option,  $a$  will not have an inverse as  $ae = aa = a$ . This contradicts the existence of an inverse axiom. You can convince yourself that the first option is consistent with all group axioms. So, there is exactly one group of size 2. In fact, as a computer scientists, you know that group very well. Do you recognize it? It is addition modulo 2—the way we are used to working with bits. This becomes more transparent if we choose to represent the group as an additive group, writing 0 for  $e$  and 1 for  $a$ . We denote this group by  $(\mathbb{Z}_2, +)$ .

One should be careful about what does it mean for two groups to be different. It is quite often the case that you're looking at this neat group you found just to realize it is a known group in disguise. That's one of the nice aspects of group theory. Just a moment ago, we figured out that our size 2 group is something we all know very well and might incorrectly consider it as being different. In the next lecture we will formalize what does it mean for two groups to be “the same”. But, for now, here is an example of two different, by all accounts, groups of size 4 (as an exercise, work out yourself all groups of size 3).

First, we have  $(\mathbb{Z}_4, +)$  - the group of addition modulo 4 over  $\{0, 1, 2, 3\}$ . Convince yourself this is indeed a group. Second, if we will “glue two independent” copies of  $(\mathbb{Z}_2, +)$  next to each other, I claim, we will be looking at a different group of size 4. Let me formalize what I mean by gluing. First, we recall the following notation - if  $A, B$  are two sets, we define  $A \times B = \{(a, b) \mid a \in A, b \in B\}$  as the product set.

**Definition 3.3** (Direct product). Let  $(G, \bullet_G)$ ,  $(H, \bullet_H)$  be two groups. We define the group  $(G \times H, \bullet)$  as follows. For every  $(g, h), (g', h')$  in  $G \times H$ , define  $(g, h) \bullet (g', h') = (g \bullet_G g', h \bullet_H h')$ .

The definition about implicitly assume that the resulting structure is a group. Convince yourself that this assertion indeed holds. One typically writes  $G \times H$  for the direct product between groups and omit the operation symbol. Note that its neutral element is  $(e_G, e_H)$  where  $e_G, e_H$  are the neutral elements of  $G, H$ , respectively. Note also that if  $G, H$  are finite, then  $|G \times H| = |G||H|$ .

With the direct product in hand, consider  $\mathbb{Z}_2 \times \mathbb{Z}_2$  which is a group of size 4. This is actually one famous group named the *Klein four-group*. There are many ways to convince yourself that this is not at all  $\mathbb{Z}_4$ . For example, every element in  $\mathbb{Z}_2 \times \mathbb{Z}_2$  is its own inverse whereas in  $\mathbb{Z}_4$ , 1 is not its own inverse (as  $1 + 1 = 2 \neq 0$ ). Turns out that there are no more groups of size 2 but to say such a thing formally we'll need to wait for the next lecture. Still, for now, here is a table of the number of groups



of a given size up to 16. We also count commutative and non-commutative groups separately to get some more insight.

Size	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Commutative	1	1	1	2	1	1	1	3	2	1	1	2	1	1	1	5
Non-commutative	0	0	0	0	0	1	0	2	0	1	0	3	0	1	0	9

By the above examples, you have probability figured out by now that for every integer  $n \geq 1$  there is a group of size  $n$ . It is the group of addition modulo  $n$  over the set  $\{0, 1, \dots, n-1\}$ . We denote this group by  $(\mathbb{Z}_n, +)$ . For every  $n$ , this is a commutative group. Observe further that according to the table, for every prime number  $p \leq 16$ ,  $(\mathbb{Z}_p, +)$  is the only group of size  $p$ . This is true in general and we will prove this soon. Looking at size 15 one can see that this is not an if and if condition

### 3.4 Some basic general properties

After seeing some groups we are ready to go abstract again and ask what kind of general properties hold for an abstract, general, group. First, we observe that a neutral element in every group is *unique*. Indeed, if  $e_1, e_2$  are both neutral then  $e_1 = e_1 e_2 = e_2$ . This allows us to say talk about *the* neutral element in a give group. An inverse of a given element is also unique. In fact, more is true. Consider a group  $G$  and a fixed element  $a \in G$ . Define the map  $\phi_a : G \rightarrow G$  by  $\phi_a(x) = ax$ . I claim that  $\phi_a$  is injective (that is, one to one). Indeed, if  $\phi_a(x) = \phi_a(y)$  then  $ax = ay$ . Let's take an inverse  $b$  of  $a$ . Then  $b(ax) = b(ay)$  and so, by the associativity law,  $(ba)x = (ba)y$ . As  $b$  is an inverse of  $a$ , we get that  $ex = ey$  and so  $x = y$ . To see that this implies the uniqueness of an inverse of a given element, consider some element  $a$ . As  $\phi_a(x)$  is injective, there is at most one element  $x$  for which  $ax = \phi_a(x) = e$  and, by the group axioms, we are guaranteed that such an element  $x$  exists. From now on we can talk about *the* inverse of an element  $a$ . If we have a multiplicative group in mind, we denote the inverse of  $a$  by  $a^{-1}$ . For an additive group we write  $-a$ . If we have a multiplicative group  $(G, \cdot)$  in mind,  $a \in G$  and  $n \in \mathbb{N}$ , you can prove that the inverse of  $a^n$  which is denoted by  $(a^n)^{-1}$ , equals to  $(a^{-1})^n$ . So we can actually extend the exponent to be an element of  $\mathbb{Z}$ . Same is true for an additive group.

### 3.5 Subgroups

In Mathematics, when you have an object that you care about that has some structure, you will typically be caring about parts of the object that also have such structure.

Wow, that was abstract. Perhaps the simplest example is the integers  $(\mathbb{Z}, +)$  and even numbers as their sub object. The structure of  $(\mathbb{Z}, +)$  is its group structure. The even numbers can also be added and they have a neutral element sitting inside. Associativity is obvious as it follows from the associativity of the integers. What we just said is that the even numbers, with respect to the addition operation is a group, which we'll denote by  $(2\mathbb{Z}, +)$ . But we really want to say that it is a *subgroup* of  $(\mathbb{Z}, +)$ . So, let's define this notion in the most natural way.

It will be useful to go back to the functional definition. We also use the following standard notation. Let  $A, B, C$  be three sets such that  $A \subseteq B$ . Let  $f: B \rightarrow C$ . We define the function  $f|_A: A \rightarrow C$  by  $f|_A(a) = f(a)$  for all  $a \in A$  and say that  $f|_A$  is the restriction of  $f$  to  $A$ .

**Definition 3.4** (Subgroup). Let  $(G, f)$  be a group. A *subgroup*  $(H, f|_{H \times H})$  of  $(G, f)$  is a group such that  $H \subseteq G$ .

Observe that as we require that  $(H, f|_{H \times H})$  is a group, we implicitly require that the image of  $f|_{H \times H}$  is contained in  $H$ .

## 3.6 Cyclic groups and generated subgroups

We want to study the properties of non-trivial groups with the simplest structure i.e we will look at groups that are generated from one element.

**Definition 3.5** (Cyclic group). A group  $G$  is *cyclic* if there exist an element  $g \in G$  that generates the whole group, i.e.  $G = \{g^i \mid i \in \mathbb{Z}\}$ , denoted  $\langle g \rangle$  or  $(g)$ .

*Remark.* Let  $g \in G$  be an element in a group  $G$  then the cyclic subgroup  $\langle g \rangle$  is the minimal subgroup of  $G$  that contains  $g$ .

*Remark.* Any cyclic group is abelian.

### Examples:

1. The trivial group with one element is cyclic and equals to  $\langle e \rangle$ .
2.  $(\mathbb{Z}, +)$  is a cyclic group of infinite size and is generated by  $\langle 1 \rangle$ .
3.  $\mathbb{Z}_n = \{0, 1, \dots, n-1\}$  with addition operation modulo  $n$  is cyclic and equals to  $\langle 1 \rangle$ .
4. The subgroup  $2\mathbb{Z}$  contains all even numbers is cyclic and equals to  $\langle 2 \rangle$ .

Naturally, one could ask to extend the definition of a subgroup that is generated by one element to a subgroup that is generated from a subset  $S$  of elements in  $G$ .

**Definition 3.6** (Subgroup). Let  $S \subseteq G$  be a subset of elements in group  $G$ . The **subgroup generated by  $S$**  is the minimal subgroup of  $G$  that contains all the elements in  $S$ , denoted by  $\langle S \rangle$ .

It follows from the definition that  $\langle S \rangle$  is formed by the intersection of all subgroups  $H < G$  that contains  $S$ :

$$\langle S \rangle = \bigcap_{S \subseteq H \leq G} H$$

**Example:**

Let  $S = \{g, h\} \subseteq G$  then the subgroup generated by  $S$  is the following:

$$\langle g, h \rangle = \{g^{i_1} h^{i_2} \dots g^{i_{n-1}} h^{i_n} \mid i_1, \dots, i_n \in \mathbb{Z}, n \in \mathbb{N}\}$$

that is,  $\langle g, h \rangle$  contains all the words that can be formed from  $g, h$  and their inverses  $g^{-1}, h^{-1}$ .

**Definition 3.7** (Element order). Let  $G$  be a group and  $g \in G$  be some element. The *order* of  $g$  defined to be the minimal integer  $0 < n$  for which  $g^n = e$ , denoted  $o(g)$ . In case  $n$  does not exist, the order of the group is infinite.

**Claim 8.** Let  $G$  be a group and  $g \in G$  be some element. The size of the cyclic group  $\langle g \rangle$  equals to the order of  $g$ , i.e.  $|\langle g \rangle| = o(g)$ .

*Proof.* Suppose the order of  $g$  is finite,  $o(g) = n$ . On the one hand,  $\{e, g, \dots, g^{n-1}\}$  contains at least  $n$  different elements: suppose to the contrary there exist  $0 \leq i < j \leq n-1$  such that  $g^i = g^j$ , then by multiplying by  $g^{-i}$  on both sides we get  $g^{j-i} = e$  but  $j-i < n$  in contradiction to  $n$  being the order of  $g$ . We conclude that  $o(g) \leq |\langle g \rangle|$ . On the other hand, we claim that  $\langle g \rangle = \{e, g, \dots, g^{n-1}\}$  is a group: every element  $g^i$  has an inverse  $g^{n-i}$  since  $g^i g^{n-i} = g^n = e$ . So this set forms a group of size  $n$  and contains  $g$  so by the definition of the cyclic group we get that  $n = o(g) \geq |\langle g \rangle|$ . All in all,  $o(g) = |\langle g \rangle|$ . □

We now want to return to the notion of cosets we saw in the last lecture, and to prove some of their properties. Let  $G$  be a group and  $H \leq G$  a subgroup, fix  $g \in G$  and we call the following a left and right coset of  $H$  in correspondence:

$$gH = \{gh \mid h \in H\}$$

$$Hg = \{hg \mid h \in H\}$$

Recall the motivation for these definitions comes from the integers  $\mathbb{Z}$  where we know that  $2\mathbb{Z}$  is a subgroup that contains all the even numbers but the odd numbers  $1 + 2\mathbb{Z}$

doesn't form a subgroup although their structures are not that far apart, thus we wish to capture these phenomena with our algebraic structures.

**Claim 9.** *Let  $H \subseteq G$ . The left cosets of the form  $gH$  for any  $g \in G$ , decompose  $G$  into disjoint sets of identical size. That is,  $\forall g_1, g_2 \in G$ :*

1. *Either  $g_1H = g_2H$  or  $g_1H \cap g_2H = \emptyset$ .*
2. *There exist a bijection  $f : g_1H \rightarrow g_2H$ .*

*Remark.* The exactly analogous claim holds for right cosets as well.

*Proof.*

1. We show that if  $g_1H \cap g_2H \neq \emptyset$  then  $g_1H = g_2H$ . Following the assumption, there exist  $h_1, h_2 \in H$  not necessarily different such that  $g_1h_1 = g_2h_2$ . Now, let  $g_1h \in g_1H$  be some element.

$$g_1h = g_1eh = g_1h_1h_1^{-1}h = (g_1h_1)(h_1^{-1}h) = (g_2h_2)(h_1^{-1}h) = g_2(h_2h_1^{-1}h) \in g_2H$$

We conclude that  $g_1H \subseteq g_2H$ , and by symmetrical arguments we would achieve  $g_2H \subseteq g_1H$ , so overall  $g_1H = g_2H$ .

2. Define  $f$  as  $x \mapsto g_2g_1^{-1}x$ .

To see its 1:1, take some  $x_1, x_2 \in g_1H$  s.t.  $g_2g_1^{-1}x_1 = g_2g_1^{-1}x_2$ , and hence

$$(g_1g_2^{-1})g_2g_1^{-1}x_1 = (g_1g_2^{-1})g_2g_1^{-1}x_2 \implies x_1 = x_2$$

The mapping is also onto: for any  $g_2h \in g_2H$  we could take  $x = g_1h \in g_1H$  so  $f(x) = g_2g_1^{-1}g_1h = g_2h$  as desired.

□

From the claim above we can conclude Lagrange Theorem.

**Theorem 3.10** (Lagrange Theorem). *Let  $G$  be a finite group and  $H \leq G$ , then  $|H|$  divides  $|G|$ .*

*Proof.* From claim 9 we conclude that  $G$  is divided into a collection of disjoint cosets of equal size, the size of each coset is  $|H|$ , so consequentially  $|H|$  divides  $|G|$ . □

**Corollary 3.11.** *Let  $G$  be a finite group such that  $|G|$  is a prime number, then  $G$  is cyclic.*

*Proof.* Take  $e \neq g \in G$ , and look at the cyclic subgroup  $\langle g \rangle \leq G$ . Due to Lagrange theorem (3.10),  $|\langle g \rangle| \mid |G| = p$ . Since  $p$  is a prime and  $\langle g \rangle$  isn't trivial then  $|\langle g \rangle| = p$ . It follows immediately that  $\langle g \rangle = G$ , from size comparison.  $\square$

**Claim 12.** *Let  $G$  be a finite group, and  $g \in G$ . Then  $g^{|G|} = e$ .*

*Proof.* Due to theorem 3.10 and claim 8, we know that  $o(g) = |\langle g \rangle|$ ,  $|\langle g \rangle| \mid |G|$ . Thus  $|G| = m \cdot o(g)$  for some  $m \in \mathbb{N}$ , and

$$g^{|G|} = g^{o(g)m} = (g^{o(g)})^m = e^m = e$$

$\square$

## 3.7 Multiplicative Group

**Definition 3.13** (Multiplicative group). For any  $n \in \mathbb{N}$ , we define  $(U_n, \cdot \text{ mod } n)$ , denoted sometimes as  $(\mathbb{Z}/n\mathbb{Z})^\times$ , as the group of all integers that coprime to  $n$ , i.e

$$U_n := \{1 \leq k \leq n \mid \gcd(k, n) = 1\}$$

**Examples:** Recall that a group is cyclic if it's generated by some element.

- $U_5 = \{1, 2, 3, 4\}$ ,  $2^2 = 4$ ,  $2^3 = 3$ ,  $2^4 = 1 \implies$  cyclic
- $U_7 = \{1, 2, 3, 4, 5, 6\}$ ,  $3^2 = 2$ ,  $3^3 = 6$ ,  $3^4 = 4$ ,  $3^5 = 5$ ,  $3^6 = 1 \implies$  cyclic
- $U_8 = \{1, 3, 5, 7\}$ ,  $3^2 = 1$ ,  $5^2 = 1$ ,  $7^2 = 1 \implies$  not cyclic

**Claim 14.** *For every  $n \in \mathbb{N}$ ,  $(U_n, \cdot \text{ mod } n)$  is indeed a group.*

*Proof.* We'll prove each property:

- Identity element:  $1 \in U_n$  is the identity for integer multiplication, and always in  $U_n$  as  $\gcd(1, n) = 1$ .
- Closure: Any  $k, k' \in U_n$  implies  $\gcd(k, n) = 1$  and  $\gcd(k', n) = 1$ , and thus also  $\gcd(kk', n) = 1$ , which in turn implies the closure property  $k \cdot k' \in U_n$ .
- Associativity: the operation is indeed associative:  $ab = ba \text{ mod } n$ .
- Inverse:  $\forall k$  define  $f_k : U_n \rightarrow U_n$  by  $x \mapsto k \cdot x$ . It's injective, as for any two inputs that holds  $f_k(x) = f_k(y)$  we got

$$k \cdot x \equiv k \cdot y \pmod{n} \implies k(x - y) \equiv 0 \pmod{n} \implies n \mid k(x - y)$$

recall that  $\gcd(n, k) = 1$ , and it follows immediately:

$$n \mid (x - y) \implies x - y \equiv 0 \pmod{n} \implies x \equiv y \pmod{n}$$

Therefore  $f_k$  is a bijection (1:1 for finite  $U_n$ ) and in particular is surjective. Hence,  $\exists k'$  s.t.  $f(k') = k \cdot k' \equiv 1 \pmod{n}$ . If we insist on finding the inverse explicitly, we can use the following theorem.

**Theorem 3.15** (Bezout identity). *Let  $a, b \in \mathbb{N}$  with  $\gcd(a, b) = d$ . Then, there exists integers  $x$  and  $y$  such that  $ax + by = d$ .*

So, in our case we get  $x$  and  $y$  such that  $kx + ny = 1$  (in particular  $x$  and  $n$  must be coprime) and this means that  $kx \equiv 1 \pmod{n}$ , i.e.  $x$  is the inverse of  $k$  in  $U_n$  and we can find it using extended Euclid's algorithm.

□

**Theorem 3.16** (Euler). *Let  $a, n \in \mathbb{N}$  such that  $\gcd(a, n) = 1$ . The Euler function  $\varphi(n)$  counts the number of positive integer that are smaller and coprime to  $n$ ,*

$$\varphi(n) = |\{1 \leq k \leq n - 1 \mid \gcd(k, n) = 1\}|$$

*Then, the following holds:*

$$a^{\varphi(n)} \equiv 1 \pmod{n}$$

*Proof.* Consider the multiplicative group  $U_n$ , and observe that  $a$  is an element in it. Following claim 12, it follows immediately that  $a^{|U_n|} = 1 \implies a^{\varphi(n)} = 1$ . □

**Corollary 3.17** (Fermat's Little Theorem). *Let  $p$  be a prime and  $a \in \mathbb{N}$ , then*

$$a^p \equiv a \pmod{p}$$

*Proof.* Note that  $\varphi(p) = p - 1$  for any prime number, and thus following 3.16,

$$a^{p-1} \equiv 1 \pmod{p} \implies a^p \equiv a \pmod{p}$$

□

### 3.8 Normal subgroups

**Definition 3.18** (Normal subgroups). Let  $G, H$  be two groups such that  $H < G$ . We say that  $H$  is a *normal subgroup* of  $G$ , denoted  $H \triangleleft G$ , if  $\forall g \in G, gH = Hg$ . In other words, any member of  $G$  induce equal right and left cosets with  $H$ .

**Definition 3.19** (Quotient group). Let  $H \triangleleft G$ . The quotient group is defined by

$$G/H := \{gH \mid g \in G\}$$

and for any two elements  $gH, g'H \in G/H$ , the operation defined by

$$gH \cdot g'H := (g \cdot g')H$$

For this to be well defined (independent of representative choice), we need:

$$\forall g, g' \in G \quad \forall h, h' \in H \quad : \quad ghH \cdot g'h'H \stackrel{?}{=} gH \cdot g'H$$

By definition  $gH \cdot g'H = gg'H$  and  $ghH \cdot g'h'H = ghg'h'H$  so equivalently we need:

$$ghg'h'H \stackrel{?}{=} gg'H$$

For normal subgroups,  $g'H = Hg'$  therefore  $\exists h'' \in H$  s.t.  $hg' = g'h''$ , hence:

$$g \underbrace{hg'}_{g'h''} h' = gg' \underbrace{h''h'}_{\in H} \in gg'H$$

Thus for normal subgroups the equality holds and the group is defined. This is called the quotient group.

As an example, observe the quotient group  $\mathbb{Z}/5\mathbb{Z}$ . Following definition,  $\mathbb{Z}/5\mathbb{Z} = \{a + 5\mathbb{Z} : a \in \mathbb{Z}\}$ . However, it is easy to see that  $a, b \in \mathbb{Z}$  such that  $a \equiv b \pmod{5}$  induce the same group, so there are essentially only 5 groups in the quotient group:

$$\mathbb{Z}/5\mathbb{Z} = \{5\mathbb{Z}, 1 + 5\mathbb{Z}, 2 + 5\mathbb{Z}, 3 + 5\mathbb{Z}, 4 + 5\mathbb{Z}\}$$

**Theorem 3.20** (Cauchy). *Let  $G$  be any finite abelian group, and  $p \in \mathbb{N}$  be some prime factor of  $|G|$ . Then there exists some  $g \in G$  such that  $o(g) = p$ .*

*Proof.* By induction on  $n = |G|$ . As  $n \geq p$  due to definition, the base case is  $n = p$ . Such a group is cyclic, so its generator has order  $p$  (in fact in this case any non-trivial element of  $G$  has order  $p$ ). Thus, we assume that for any group with size smaller than  $n$  the theorem holds, and prove it for  $n$ .

Take some  $x \in G$ , and denote by  $X = \langle x \rangle$  the cyclic group generated by  $x$ . If  $p \mid |X|$ , then since  $x^{|X|} = e$  (by definition), it follows immediately that  $x^{|X|/p}$  has order  $p$ .

Now, we deal with the case  $p \nmid |X|$ .  $G$  is abelian, so  $X$  is normal. Recall that  $|G| = |G/X| \cdot |X|$ , so we must have  $p \mid |G/X|$ . Consider the quotient group  $G/X$ : it is abelian, and since  $|G/X| < n$  the induction hypothesis apply on it. That is, there exist  $g \in G$  such that  $gX$  has order  $p$  at the quotient group. Observe that for  $o(g)$ , the order of  $g$  at  $G$ ,

$$(gX)^{o(g)} = g^{o(g)}X = eX = X$$

We must have  $p \mid o(g)$ , otherwise we can write  $o(g) = qp + r$  with  $q, r \in \mathbb{Z}$  and  $1 \leq r < p$ , which means that

$$X = (gX)^{o(g)} = (gX)^{qp+r} = (gX)^{qp}(gX)^r = ((gX)^p)^q(gX)^r = X(gX)^r = (gX)^r$$

which is a contradiction to  $p$  being the order of  $gX$  in  $G/X$ . Finally,  $g^{o(g)/p}$  has order  $p$ , concluding the proof.  $\square$

### 3.9 Homomorphisms

**Definition 3.21** (Homomorphism). Let  $G$  and  $H$  be two groups, and denote their operations as  $\cdot_G, \cdot_H$  respectively. A map  $\varphi : G \rightarrow H$ , is called *homomorphism*, if it preserves the structure of the groups, under their operations:

$$\forall a, b \in G : \varphi(a \cdot_G b) = \varphi(a) \cdot_H \varphi(b)$$

An injective homomorphism, is called *monomorphism*.

A surjective homomorphism, is called *epimorphism*.

A bijective homomorphism, is called *isomorphism*.

For example,

1.  $\varphi : \mathbb{Z} \rightarrow \mathbb{Z}_5$  defined by  $\varphi(x) = x \pmod{5}$ .
2.  $\psi : (\mathbb{R}, +) \rightarrow (\mathbb{R} \setminus \{0\}, \cdot)$  defined by  $\psi(z) = e^z$ .

Let  $G$  and  $H$  be two groups and  $\varphi : G \rightarrow H$  be an homomorphism. Observe the following properties:

1.  $\varphi(e) = e$ .
2.  $\varphi(g^{-1}) = \varphi(g)^{-1}$ .



$$3. \varphi(G) = \{\varphi(g) : g \in G\} < H.$$

**Definition 3.22** (Group Isomorphism). Let  $G$  and  $H$  be two groups. If there exists an isomorphism  $\varphi : G \rightarrow H$ , then we say that  $G$  is *isomorphic* to  $H$ , denoted  $G \cong H$ .

Observe that  $\cong$  is an equivalence relation over the groups, this means that

1. **Reflexive.** If  $G$  is a group, then  $G \cong G$ : simply take  $\varphi$  to be the identity  $\varphi(a) = a$ .
2. **Symmetric.** If  $G$  and  $H$  are groups such that  $G \cong H$ , then  $H \cong G$ : since  $\varphi$  is bijective,  $\varphi^{-1}$  is the inverse isomorphism.
3. **Transitive.** If  $G$ ,  $H$  and  $K$  are groups such that  $G \cong H$  and  $H \cong K$ , then  $G \cong K$ : simple composition of isomorphisms.

For example, consider the multiplicative group  $U_5 = \{1, 2, 3, 4\}$ . Observe that  $2^2 = 4, 2^3 = 3$ . This means that 2 is a generator of  $U_5$ , i.e.  $U_5 = \langle 2 \rangle$  is a cyclic group. Now, consider the additive group  $\mathbb{Z}_4 = \{0, 1, 2, 3\}$ , then  $U_5 \cong \mathbb{Z}_4$  using the mapping  $\varphi : U_5 \rightarrow \mathbb{Z}_4$  defined by

$$\varphi(1) = 0, \varphi(2) = 1, \varphi(3) = 3, \varphi(4) = 2$$

Alternatively, we can take the multiplicative group  $U_8 = \{1, 3, 5, 7\}$ . This group is not cyclic since the order of each non-trivial element is 2. It is possible to show that this group is isomorphic to the product group  $\mathbb{Z}_2 \times \mathbb{Z}_2$ . We can look at the elements of  $U_8$  in a binary representation using only 3 bits:

$$1 = 001, 3 = 011, 5 = 101, 7 = 111$$

Observe that all the numbers have 1 as their least significant bit (they are all odd), so we can map each number to the first two bits

$$\varphi(1) = (0, 0), \varphi(3) = (0, 1), \varphi(5) = (1, 0), \varphi(7) = (1, 1)$$

The mapping  $\varphi : U_8 \rightarrow \mathbb{Z}_2 \times \mathbb{Z}_2$  is clearly bijective, you can manually check that  $\varphi$  is an homomorphism.

Another interesting example arises when looking at a group  $G$  with a normal subgroup  $N \triangleleft G$ . We get an epimorphism  $\varphi : G \rightarrow G/N$  by simply defining

$$\forall g \in G : \varphi(g) = gN$$

**Definition 3.23** (Kernel). Let  $G$  and  $H$  be two groups, and  $\varphi : G \rightarrow H$  be an homomorphism. The *kernel* of  $\varphi$  is defined to be all the elements of  $G$  which goes to  $e_H$ .

$$\text{Ker}(\varphi) = \varphi^{-1}(\{e_H\}) = \{g \in G : \varphi(g) = e_H\}$$

**Lemma 3.24.** Let  $G$  and  $H$  be two groups, and  $\varphi : G \rightarrow H$  be an homomorphism, then

$$\varphi \text{ is a monomorphism} \iff \text{Ker}(\varphi) = \{e\}$$

*Proof.* Suppose that  $\text{Ker}(\varphi) = \{e\}$ . Let  $a, b \in G$  such that  $\varphi(a) = \varphi(b)$ . Then

$$\varphi(ab^{-1}) = \varphi(a) \cdot \varphi(b^{-1}) = \varphi(a) \cdot \varphi(b)^{-1} = \varphi(a) \cdot \varphi(a)^{-1} = e$$

this means that  $ab^{-1} \in \text{Ker}(\varphi)$ , i.e.  $ab^{-1} = e$  and finally  $a = b$ . This means that  $\varphi$  is injective and thus a monomorphism. The other direction is trivial.  $\square$

# LECTURE 4

## QUOTIENT GROUP

---

### 4.1 Group Theory - Continued

#### 4.1.1 From previous lecture

**Claim 1.** Let  $H < G$ .  $H \triangleleft G \iff \exists \varphi$  homomorphism such that  $\text{Ker}\varphi = H$

*Proof.*

$\Leftarrow$

We proved every kernel is normal, thus if  $H = \text{Ker}\varphi$  then  $H = \text{Ker}\varphi \triangleleft G$ .

$\Rightarrow$

Consider the natural homomorphism  $\varphi : G \rightarrow G/H$  defined by  $\varphi(g) = gH$ . Then:

$$g \in \text{Ker}\varphi \iff \varphi(g) = H \iff gH = H \iff g \in H \quad (1)$$

Hence  $\text{Ker}\varphi = H$ , as required.  $\square$

#### 4.1.2 First Isomorphism Theorem

**Theorem 4.2** (First Isomorphism Theorem). Let  $\varphi : G \rightarrow \text{Im}\varphi$  be a homomorphism. Then  $G/\text{Ker}\varphi \simeq \text{Im}\varphi$

#### Example:

Consider  $\varphi : \mathbb{Z} \rightarrow \mathbb{Z}_5$  defined by  $\varphi(x) = x \bmod 5$ . Then  $\text{Ker}\varphi \simeq 5\mathbb{Z}$ . Therefore  $\mathbb{Z}/5\mathbb{Z} \simeq \mathbb{Z}_5$

*Proof.* We will find an isomorphism from  $G/\text{Ker}\varphi$  to  $\text{Im}\varphi$ .

Denote  $K = \text{Ker}\varphi$  and define  $\psi(gK) = \varphi(g)$ . We will prove this function is well defined and that it is an isomorphism.

Let  $g_1, g_2$  such that  $g_1K = g_2K$ . Then  $g_2 = g_1k$  for some  $k \in K$ . Hence  $\psi(g_2K) = \psi(g_1kK) = \varphi(g_1k) = \varphi(g_1)\varphi(k) = \varphi(g_1) = \psi(g_1K)$ , where the one-before-last equality follows from that  $k \in K$ .

$\psi$  is a homomorphism:  $\psi(gK \cdot hK) = \psi(ghK) = \varphi(gh) = \varphi(g) \cdot \varphi(h) = \psi(gK) \cdot \psi(hK)$ . We used the fact that  $\varphi$  is a homomorphism itself.

$\psi$  is onto: Let  $x \in \text{Im}\varphi$ . Then there is a  $g$  such that  $\varphi(g) = x$ . Therefore  $\psi(gK) = \varphi(g) = x$  as required.

$\psi$  is one-to-one: We will prove the kernel of  $\psi$  is trivial.  $gK \in \text{Ker}\psi \iff \psi(gK) = e \iff \varphi(g) = e \iff g \in \text{Ker}\varphi = K \iff gK = K$ .

We proved  $\psi$  is a group isomorphism as required. □

**Theorem 4.3** (Cyclic Group Classification). *Let  $G$  be a cyclic group. Then  $G$  is isomorphic to either  $\mathbb{Z}$  or  $\mathbb{Z}_n$*

*Proof.* As  $G$  is cyclic, every element in  $G$  can be expressed as a power of some  $g$ .

We define  $\varphi : \mathbb{Z} \rightarrow \langle g \rangle$  as follow:  $n \mapsto g^n$ .  $\varphi$  is a homomorphism as  $\varphi(n + m) = g^{n+m} = g^n \cdot g^m = \varphi(n) \cdot \varphi(m)$ .

$\varphi$  is onto by the definition of the cyclic group. Therefore, by the first isomorphism theorem we deduce that  $\mathbb{Z}/\text{Ker}\varphi \simeq G$ .

Note that  $H = \text{Ker}\varphi$  is a subgroup of  $\mathbb{Z}$ . If  $H = \{0\}$ , then  $G \simeq \mathbb{Z}/\{0\} \simeq \mathbb{Z}$  as required.

Otherwise, let  $n$  be the smallest positive number in  $H$  and let  $m$  be some integer in  $H$ . We can divide  $m$  by  $n$  with a remainder:  $m = nq + r$  for some  $0 \leq r < n$ . Thus  $r = m - nq \in H$  but  $r < n$ , then  $r = 0$  due to the minimality of  $n$ . Then  $m$  is a multiple of  $n$ , and  $H = n\mathbb{Z}$ .

We conclude that in this case  $G = \mathbb{Z}/n\mathbb{Z} \simeq \mathbb{Z}_n$  □

**Theorem 4.4** (The Correspondence Theorem). *Let  $N \triangleleft G$ . There is a onto and one-to-one correspondence between subgroups of  $G/N$  to subgroups of  $G$  containing  $N$ .*

Example:

**Theorem 4.5** (Abelian Groups Classification). *Let  $G$  be a finite abelian group. Then there exist  $g_1, \dots, g_n$  such that  $G \simeq \langle g_1 \rangle \times \dots \times \langle g_n \rangle$  (No Proof)*

## 4.2 Ring Theory

**Definition 4.6** (Ring). We say that  $(R, +, \cdot)$  is a ring if:

1.  $(R, +)$  is an abelian group, and  $0$  is its neutral element
2.  $a, b \in R \Rightarrow a \cdot b \in R$
3.  $\forall a, b, c \in R : (a \cdot b) \cdot c = a \cdot (b \cdot c)$
4.  $\exists 1 \in R \forall a \in R : 1 \cdot a = a \cdot 1$

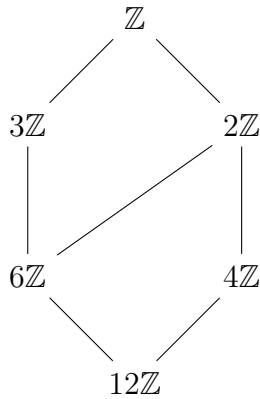


Figure 3: Subgroups of  $\mathbb{Z}$  containing  $12\mathbb{Z}$

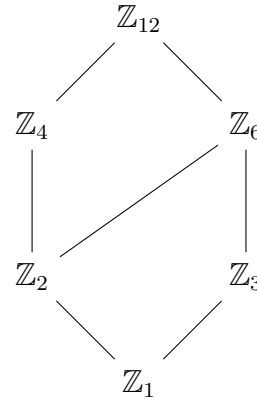


Figure 4: Subgroups of  $\mathbb{Z}/12\mathbb{Z} \simeq \mathbb{Z}_{12}$

5.  $\forall a, b, c \in R : a \cdot (b + c) = a \cdot b + a \cdot c$
6.  $\forall a, b, c \in R : (a + b) \cdot c = a \cdot c + b \cdot c$
7.  $0 \neq 1$

**Definition 4.7** (Commutative Ring). Let  $(R, +, \cdot)$  be a ring. We say that  $(R, +, \cdot)$  is a commutative ring if  $\forall a, b \in R \ a \cdot b = b \cdot a$ .

*Remark.* From now on we will talk about commutative rings (CR).

*Example 4.8.*  $(\mathbb{Z}, +, \cdot)$  is a CR. note that  $2 \cdot 3 = 0$ .

**Definition 4.9** (Zero Divisor). We say that  $r \in R \setminus \{0\}$  is a zero divisor if there exists  $s \in R \setminus \{0\}$  such that  $r \cdot s = 0$

**Definition 4.10** (Integral Domain). Let  $R$  be a CR. We say that  $R$  is an integral domain if  $\forall r \in R, \ r$  is not a zero divisor.

**Definition 4.11** (Field). Let  $R$  be a CR. We say that  $R$  is a field, if  $\forall r \in R \ \exists s \in R : r \cdot s = 1$

*Example 4.12.*

**Corollary 4.13.** *Let  $F$  be a field. Then,  $F$  is an integral domain.*

*Proof.* Let  $x, y \in F$  such that  $xy = 0$ . If  $x \neq 0$  there exist  $x^{-1} \in R$  such that  $x^{-1}x = 1$ . Hence,  $y = x^{-1}xy = x^{-1}0 = 0$ . Otherwise,  $x = 0$ . Therefore, there aren't any zero divisors in  $F$  which implies that it is an integral domain.  $\square$

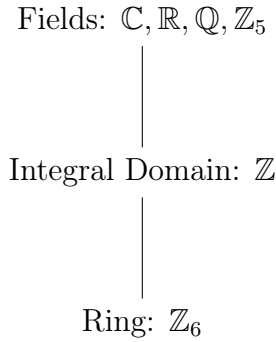


Figure 5: Hierarchy of fields, integral domains, and rings.

### 4.3 Properties of Commutative Rings

1.  $0 \cdot a = a \cdot 0 = 0$
2.  $a \cdot (-b) = (-a) \cdot b = -ab$
3.  $(-a)(-b) = ab$
4.  $(-1)a = -a$
5.  $(a + b)^2 = a^2 + 2ab + b^2$

*Proof.*

1. Note that  $0 \cdot a = (0 + 0) \cdot a = 0 \cdot a + 0 \cdot a$  using distribution. Therefore  $0 \cdot a = 0$ , and by commutativity,  $a \cdot 0 = 0$  as well.
2. To prove that  $a \cdot (-b) = -ab$ , we need to prove that  $a \cdot (-b) + ab = 0$ . By distribution,  $a \cdot (-b) + ab = a(-b + b) = a \cdot 0 = 0$  by property 1.
3.  $(-a)(-b) = -(-a)b = -(-ab)$  by repeated application of property 2. As the negative of a negative is the original number by definition,  $-(-ab) = ab$  as required.
4. By property 2,  $(-1) \cdot a = -(1 \cdot a) = -a$ .
5. Expanding the expression we get  $(a + b)^2 = a^2 + ab + ba + b^2$ , and as the ring is commutative, this is equal to  $a^2 + 2ab + b^2$

□

*Remark.*  $\mathbb{Z}_p$  is a field for prime  $p$

*Proof.* We've shown that  $\mathbb{Z}_p$  is a commutative ring. It remains to show that each element  $\neq 0$  has a multiplicative inverse.

Indeed, by Fermat's little theorem, for each  $0 \neq x \in \mathbb{Z}_p$  it holds that  $x^{p-1} = 1 \pmod p$ . Therefore  $x^{p-2}$  is the multiplicative inverse of  $x$ .  $\square$

**Definition 4.14** (Alternative definition of Field). A tuple  $(R, +, \cdot)$  is a field if:

1.  $(R, +)$  is an abelian group (with 0 as the identity element)
2.  $(R \setminus \{0\}, \cdot)$  is an abelian group (with 1 as the identity element)
3.  $(a + b)c = ac + bc$

**Definition 4.15** (Sub-Ring). A subset  $S$  of a ring  $(R, +, \cdot)$  is a subring if it is a ring under  $(+, \cdot)$

**Definition 4.16** (Ring Homomorphism). Let  $R, R'$  be rings. A function  $\varphi : R \rightarrow R'$  is a ring homomorphism if:

1.  $\forall r, s \in R, \varphi(r + s) = \varphi(r) + \varphi(s)$
2.  $\forall r, s \in R, \varphi(r \cdot s) = \varphi(r) \cdot \varphi(s)$
3.  $\varphi(1) = 1$

**Claim 17.** Let  $\varphi : R \rightarrow R'$  be a ring homomorphism.

1.  $\varphi$  is one-to-one  $\iff \text{Ker}\varphi = \{0\}$
2.  $\forall k \in \text{Ker}\varphi, \forall r \in R, \varphi(k \cdot r) = 0$

*Proof.*

1. Assume  $\varphi$  is one-to-one, and let  $r \in \text{Ker}\varphi$ . By definition,  $\varphi(r) = 0$ , but  $\varphi(0) = 0$  as  $\varphi$  is an homomorphism. Since  $\varphi$  is one-to-one,  $r = 0$ .

Now assume  $\text{Ker}\varphi = \{0\}$  and let  $x, y \in R$  such that  $\varphi(x) = \varphi(y)$ . Therefore  $\varphi(x - y) = \varphi(x) - \varphi(y) = 0$  since  $\varphi$  is an homomorphism. We deduce that  $x - y \in \text{Ker}\varphi$ , thus  $x - y = 0$  or  $x = y$  as required.

2.  $\varphi(k \cdot r) = \varphi(k)\varphi(r) = 0 \cdot \varphi(r) = 0$

$\square$

**Definition 4.18** (Ring Isomorphism). Let  $R, R'$  be rings. We say the rings are isomorphic and denote  $R \simeq R'$  if there is an isomorphism  $\varphi : R \rightarrow R'$

**Definition 4.19** (Ideal).  $I \subset R$  for some ring  $R$  is an ideal if

1.  $I$  is an additive subgroup
2.  $\forall i \in I, r \in R \rightarrow ir \in I$

Two trivial ideals always exist:  $R, \{0\}$ . Another example is that for  $R = (\mathbb{Z}, +, \cdot)$ .

**Definition 4.20** (Principal Ideal). An ideal  $\langle r \rangle$  of  $R$  will be called a *Principal Ideal* if it is generated by a single element  $r$  through multiplication by every element of  $R$ .

$$\langle r \rangle = \{rs \mid s \in R\} = rR$$

For a set  $S \subset R$  we mark

$$\langle S \rangle = \left\{ \sum s_i r_i \mid r_i \in R \right\}$$

and for  $s, t \in R$  we mark

$$\langle s, t \rangle = sR + tR$$

**Definition 4.21** (Principal Ideal Domain (PID)). An integral domain where every ideal is principal will be called a *Principal Ideal Domain*.

**Definition 4.22** (Quotient Ring). Let  $(R, +, \cdot)$  be a ring and  $I \subset R$  and ideal. We set

$$R/I = \{r + I \mid r \in R\}$$

and we define the  $\cdot$  operation by the representatives' operation:

$$(r_1 + I) \cdot (r_2 + I) = r_1 \cdot r_2 + I$$

(because we want the "quotient group" to also represent the multiplication)

**Theorem 4.23** (The First Isomorphism). *Given an epimorphism  $\varphi : R \rightarrow Im(\varphi)$ ,*

$$Im(\varphi) \simeq R/Ker(\varphi)$$

For a given ideal  $I$  The *Correspondence Theorem* for rings will provide a two way correspondence between ideals of  $R$  that contain  $I$  and ideals of  $R/I$ .

### 4.3.1 Gaussian Integers

$$\mathbb{Z}[i] = \{a + bi \mid a, b \in \mathbb{Z}\}, i^2 = -1$$



We feel quite comfortable saying that  $\mathbb{Z}[i]$  is an integral domain. Note that 5 is not prime:

$$5 = (2 + i)(2 - i)$$

**Definition 4.24** (Maximal Ideal). Let  $R$  be a commutative ring. An ideal  $I \neq R$  is *maximal* if the ideals containing  $I$  are  $I$  and  $R$ .

For example,  $6\mathbb{Z} \subset \frac{2\mathbb{Z}}{3\mathbb{Z}} \subset \mathbb{Z}$ . So  $2\mathbb{Z}$  and  $3\mathbb{Z}$  are maximal ideals while  $6\mathbb{Z}$  is not. More generally, for  $p|n$ ,  $n\mathbb{Z} \subset p\mathbb{Z} \subset \mathbb{Z}$ .

*Remark.* The maximal ideals of  $\mathbb{Z}$  are  $\{p\mathbb{Z} \mid p \text{ is prime}\}$ .  $p\mathbb{Z}$  is clearly maximal because a  $p\mathbb{Z} \subset n\mathbb{Z}$  will give  $n \mid p$  in contrast. W.L.O.G. a maximal ideal in  $\mathbb{Z}$  is of the form  $n\mathbb{Z}$  (prime ideal) because otherwise we could split into  $m\mathbb{Z}, k\mathbb{Z} \subset n\mathbb{Z}$  (note: this is trivial from the fact that  $\mathbb{Z}$  is a PID). If  $n\mathbb{Z}$  is an ideal and  $n$  is not prime, then for  $p|n$ ,  $n\mathbb{Z} \subset p\mathbb{Z}$ , but that means its not maximal, so  $n$  must be prime.

*Remark.* We can say the abstraction of *prime numbers* for rings are *maximal ideals*.

### 4.3.2 The orthogonality of addition and multiplication

Let  $R$  be a commutative ring and  $A \subseteq R$ . We define

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

$$A \cdot A = \{a \cdot b \mid a, b \in A\}$$

For example, for  $A = \{0, 1, 2, \dots, n - 1\}$ ,

$$A + A = \{0, \dots, 2n - 2\}$$

And for  $B = \{2^0, 2^1, 2^2, \dots, 2^{n-1}\}$ ,

$$B \cdot B = \{2^0, 2^1, \dots, 2^{2n-2}\}$$

We notice that  $|A| \leq |A + A|$  and  $|A \cdot A| \leq |A|^2$ . Is there some  $A$  such that  $|A| = n$ , and both  $|A + A|$  and  $|A \cdot A|$  and not "large"? Perhaps even  $\max(|A + A|, |A \cdot A|) < |A| \cdot \log |A|$ ? Apparently *no*. Actually:

$$\forall A, \max(|A + A|, |A \cdot A|) \geq |A|^{1+\epsilon}$$

## LECTURE 5

# GROUP HOMOMORPHISM

---

In the previous lecture we've started learning ring theory. We've defined (commutative) rings, integral domains, fields and idelas. We've also seen that the notion of ideal in rings is similar to that of normal subgroups in group theory. At the end of the lecture, we've seen the first homomorphism theorem for rings and the correspondence theorem for rings.

In this lecture we will continue to investigate ring theory.

### 5.1 Fields and PID

In the following section we will show that every field is a PID. In fact, we will show even a stronger result, that every field has only the trivial ideals, that is, the ideals  $\langle 0 \rangle = \{0\}$  and  $\langle 1 \rangle = R$ .

**Claim 1.** *A commutative ring  $R$  is a field if and only if its only ideals are  $\langle 0 \rangle$  and  $\langle 1 \rangle$ .*

*Proof.* Assume that  $R$  is a field and let  $I \subseteq R$  be an ideal. If  $I = \{0\}$  we're done. Otherwise there exists  $x \in I$  such that  $x \neq 0$ . Let  $x^{-1} \in F$  be the inverse of  $x$ . Then, since  $I$  is an ideal,  $x \cdot x^{-1} = 1 \in I$ . It follows that  $I = \langle 1 \rangle$ , as required.

Assume that  $R$  is a commutative ring whose only ideals are  $\langle 0 \rangle$  and  $\langle 1 \rangle$ . Let  $x \in R$  be such that  $x \neq 0$ , and look at  $\langle x \rangle$ . Then necessarily  $\langle x \rangle = R$ , so  $1 \in \langle x \rangle$  and  $1 = x \cdot y$  for some  $y \in R$ . It follows that every  $0 \neq x \in R$  has an inverse so  $R$  is a field.  $\square$

**Corollary 5.2.** *Every field is a PID.*

### 5.2 Primes and Irreducibles

**Motivation.** How would you define a prime number in  $\mathbb{Z}$ ?

- **1st attempt:**  $p$  is a prime if and only if, for any  $a, b \in \mathbb{Z}$ , when  $p = ab$  then either  $a$  or  $b$  are  $\pm 1$ .
- **2nd attempt:**  $p$  is a prime if and only if, for any  $a, b \in \mathbb{Z}$ , if  $p|ab$  then either  $p|a$  or  $p|b$ .

It turns out that in  $\mathbb{Z}$  both definitions are equal, but we shall see that in other number system this does not necessarily true. Indeed the first attempt defines *irreducible* numbers while the second defines *primes*, and we will see that in integral domains being a prime implies irreducibility.

**Definition 5.3.** Let  $R$  be an integral domain. For any  $r, s \in R$ , we say that  $r|s$  ( $r$  divides  $s$ ) if and only if there exists  $t \in R$  such that  $s = rt$ .

*Remark.* Note that  $r|s$  if and only if  $\langle s \rangle \subseteq \langle r \rangle$ .

**Definition 5.4.** Let  $R$  be an integral domain. An element  $r \in R$  is a *unit* if there exists  $s \in R$  such that  $rs = 1$ .

*Remark.* Note that  $r$  is a unit if and only if  $r|1$ . This happens if and only if  $\langle r \rangle = \langle 1 \rangle$ .

**Example.** In a field, all non-zero elements are units (since they have an inverse). This is a sufficient condition for a commutative  $R$  ring to be a field, that is, if every element is unit then any element has an inverse, so  $R$  is a field.

**Another example.** For an integral domain  $R$  define  $U(R) = \{\text{units in } R\}$ . Then  $U(R)$  is a group under multiplication, since the multiplication of two units is a unit.

**Definition 5.5.** Let  $R$  be an integral domain. A non-zero and non-unit element  $m \in R$  is called *irreducible* if  $m = ab$  implies that  $a$  or  $b$  is a unit.

**Definition 5.6.** Let  $R$  be an integral domain. A non-zero and non-unit element  $p \in R$  is called *prime* if for any  $a, b \in R$ ,  $p|ab$  implies that  $p|a$  or  $p|b$ .

**Lemma 5.7.** *Let  $R$  be an integral domain, then every prime is irreducible.*

*Proof.* Let  $p \in R$  is a prime and assume that  $p = ab$ , so  $p|ab$ . Since  $p$  is prime we can assume wlog that  $p|a$ . Thus, there exists  $c$  such that  $a = pc$ , thus  $p = pcb$ , so  $p(1 - cb) = 0$  and since  $R$  is integral domain then  $p = 0$  or  $1 - cb = 0$ , so  $cb = 1$ . Since  $p \neq 0$  it follows that  $b$  is a unit.  $\square$

### 5.3 Constructing ideals from existing ones

Let  $R$  be a commutative ring and  $I, J$  are ideals of  $R$ , the following compositions of  $I, J$  are also ideals:

- $I \cap J$

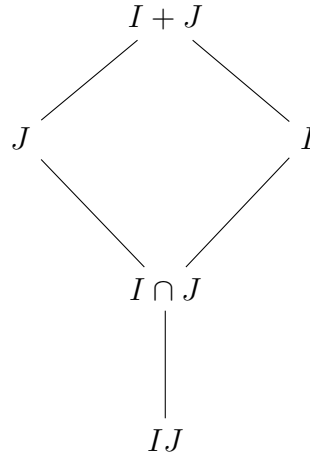


Figure 6: Relations between composition of ideals

- $IJ = \{\sum_{t=1}^n i_t j_t \mid n \in \mathbb{N}, i_1, \dots, i_n \in I, j_1, \dots, j_n \in J\}$
- $I + J = \{i + j \mid i \in I, j \in J\}$  (smallest ideal that contains  $I, J$ )

Figure 6 present graphically the relations between composition of ideals.

**Example.** Let  $R = \mathbb{Z}$ ,  $I = n\mathbb{Z}$  and  $J = m\mathbb{Z}$ . In this case we have  $I + J = \langle \gcd(n, m) \rangle$ ,  $I \cap J = \langle \text{lcm}(n, m) \rangle$  and  $IJ = \langle nm \rangle$ . For example, if  $I = 6\mathbb{Z}$  and  $J = 4\mathbb{Z}$  then  $I + J = \langle 2 \rangle$ ,  $I \cap J = \langle 12 \rangle$  and  $IJ = \langle 24 \rangle$ .

Sometimes we would like to look at an ideal that is generated by two elements of a commutative ring  $R$ . In this case, for  $a, b \in R$  we have

$$\langle a, b \rangle \triangleq aR + bR = \langle a \rangle + \langle b \rangle.$$

This is the smallest ideal with respect to inclusion that contains both  $a$  and  $b$ . Furthermore, for two elements  $a, b \in R$  we have  $\langle a \rangle \langle b \rangle = \langle ab \rangle$ .

*Proof Sketch.* Note that any element in  $\langle a \rangle \langle b \rangle$  is a sum of elements the form  $ar \cdot bs$  for  $r, s \in R$ . Each of these elements is in  $(ab)R$ , since  $ar \cdot bs = (ab)rs$ . Hence the sum of those elements is a sum of elements in  $\langle ab \rangle$  and every such sum is in  $\langle ab \rangle$  since any ideal is closed under addition. The other direction is similar.  $\square$

**Example.** Recall that  $\mathbb{Z}[x]$  is a commutative ring that includes all polynomial with integer coefficients. Let  $I \subseteq \mathbb{Z}[x]$  be the following ideal

$$I = \{\text{polynomials with an even constant term}\}.$$

One can show that  $I$  is generated by the polynomials  $x$  and  $2$ , that is  $I = \langle x, 2 \rangle = x\mathbb{Z}[x] + 2\mathbb{Z}[x]$ .

Recall that for two ideals  $J$  and  $K$ , we defined  $JK$  as all finite sums of terms of the form  $jk$  for  $j \in J$  and  $k \in K$ . The ideal  $I$  that we've just defined is a good example of why the finite sums are needed. Look at  $I^2$ , and note that  $x^2 \in I^2$  and  $4 \in I^2$ . Since  $I^2$  should also be an ideal, we would like to have  $x^2 + 4 \in I^2$ . Note that there are no two polynomials  $p, q \in I$  such that  $p(x) \cdot q(x) = x^2 + 4$ , so without taking the finite sums, this element wouldn't be in  $I^2$  and we wouldn't get an ideal.

*Fact 5.8.* Let  $R$  be a commutative ring and let  $I, J$  and  $K$  be ideals.

1.  $I(J + K) = IJ + IK$ .
2.  $I + J = R$  implies  $I \cap J = IJ$ .

## 5.4 Maximal Ideals and Prime Ideals

In the following we present the notions of maximal ideals and prime ideals.

**Definition 5.9.** Let  $R$  be a commutative ring. An ideal  $M \neq R$  is *maximal* if for any ideal  $N$  such that  $M \subseteq N \subseteq R$  then either  $N = M$  or  $N = R$ .

**Definition 5.10.** Let  $R$  be a commutative ring. An ideal  $P \neq R$  is *prime* if for any  $a, b \in R$ ,  $ab \in P$  implies  $a \in P$  or  $b \in P$ .

The following claim shows the connection between prime element of a ring and prime ideals.

**Claim 11.** *Let  $R$  be a commutative ring. A non-zero element  $p \in R$  is prime if and only if  $\langle p \rangle$  is prime.*

*Proof.* Assume that  $p$  is prime. Assume that  $ab \in \langle p \rangle$ . Recall that  $ab \in \langle p \rangle$  if and only if  $p|ab$ . Since  $p$  is prime, it holds that either  $p|a$  or  $p|b$ . Now,  $p|a$  if and only if  $a \in \langle p \rangle$  and  $p|b$  if and only if  $b \in \langle p \rangle$ , so either  $a \in \langle p \rangle$  or  $b \in \langle p \rangle$ .

Assume that  $\langle p \rangle$  is prime. Assume that  $p|ab$  for some  $a, b \in R$ . It follows that  $ab \in \langle p \rangle$ . Hence  $a \in \langle p \rangle$  or  $b \in \langle p \rangle$ , then either  $p|a$  or  $p|b$ . □

**Theorem 5.12.** *Let  $R$  be a commutative and let  $M$  be an ideal. Then  $M$  is maximal if and only if  $R/M$  is a field.*

*Proof Sketch.* Recall that  $R/M$  is a field if and only if  $R/M$  has only trivial ideals. From the correspondence theorem, every ideal of  $R/M$  has one-to-one correspondence to the ideals of  $R$  containing  $M$ . Hence if  $M$  is maximal then  $R/M$  has only two ideals, which are the trivials, and if  $R/M$  is a field then  $M$  is contained in exactly two ideals- itself and  $R$ .  $\square$

**Theorem 5.13.** *Let  $R$  be a commutative ring and let  $P \subseteq R$  be an ideal. Then  $P$  is prime if and only if  $R/P$  is integral domain.*

The last two theorems implies the following.

**Corollary 5.14.** *Let  $R$  be a commutative ring. Then every maximal ideal is prime.*

*Proof.* Let  $P$  be a maximal ideal. Then  $R/P$  is a field, hence integral domain, so  $P$  is a prime.  $\square$

An interesting fact is that the other direction is also true in PIDs, as the following claim implies and as we'll see in HW.

**Claim 15.** *Let  $R$  be an integral domain and let  $0 \neq r \in R$ . Then  $r$  is irreducible if and only if  $\langle r \rangle$  is maximal w.r.t. principal ideals.*

*Proof.* We will show only that if  $\langle r \rangle$  is maximal w.r.t. principal ideals then  $r$  is irreducible. The other direction appears as an exercise.

Assume that  $\langle r \rangle$  is maximal w.r.t. principal ideals. Let  $r = ab$ . Then, since  $a|r$ , it holds that  $\langle r \rangle \subseteq \langle a \rangle$ , and since  $\langle r \rangle$  is maximal w.r.t. principal ideals it must be the case that  $\langle a \rangle = \langle r \rangle$  or  $\langle a \rangle = R$ . If  $\langle a \rangle = R$  then  $a$  is a unit. If  $\langle a \rangle = \langle r \rangle$  then  $a = rc$  so  $r = rc b$  and  $r(1 - cb) = 0$ . Since  $R$  is an integral domain, either  $r = 0$  or  $1 - cb = 0$ . Since  $r \neq 0$  it follows that  $cb = 1$  so  $b$  is a unit.  $\square$

Figure 7 summarizes the relations between prime and maximal ideals of an integral domain. Note that since in principal ideal domain every ideal is principal, it follows that if  $r$  is irreducible then  $\langle r \rangle$  is maximal.

## 5.5 Gaussian Integers

In the following we define a commutative ring, and show some of its interesting properties. Define  $\mathbb{Z}[i] = \{a + bi \mid a, b \in \mathbb{Z}\}$  where  $i^2 = -1$ . Note that this is a commutative ring.

We show that some irreducible elements of  $\mathbb{Z}$  are reducible in  $\mathbb{Z}[i]$ . We know that 5 is irreducible in  $\mathbb{Z}$  but in  $\mathbb{Z}[i]$  it holds that  $5 = (2 + i)(2 - i)$ . Note that  $(2 + i)$  and

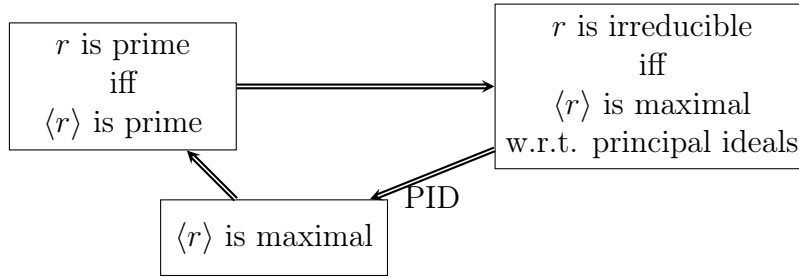


Figure 7: The relations between prime and maximal ideals of an integral domain.

$(2 - i)$  are not units, so 5 is reducible. Indeed, assume that there exists  $a, b \in \mathbb{Z}$  such that  $1 = (a + bi)(2 + i)$ . Then

$$1 = (a + bi)(2 + i) = (2a - b) + i(a + 2b),$$

and we get  $2a - b = 1$  and  $a + 2b = 0$  which has no solution in  $\mathbb{Z}$ . A similar argument shows that  $2 - i$  is not unit. This is an interesting result, by adding elements to  $\mathbb{Z}$  we made an irreducible element reducible.

We show that some irreducible elements of  $\mathbb{Z}$  remain irreducible. In particular, we show that 3 is irreducible in  $\mathbb{Z}[i]$ . We even show a stronger claim, that  $\langle 3 \rangle$  is maximal. Assume that  $\langle 3 \rangle \subsetneq I \subseteq \mathbb{Z}[i]$  where  $I$  is an ideal. Then there exists some  $r + si \in I$  such that  $r + si \notin \langle 3 \rangle$ . Since  $\langle 3 \rangle = \{3a + 3bi \mid a, b \in \mathbb{Z}\}$ , it holds that either  $r$  or  $s$  are not divisible by 3.

Look at  $t = r^2 + s^2$ , and notice that  $0^2 = 0 \pmod{3}$ ,  $1^2 = 1 \pmod{3}$  and  $2^2 = 1 \pmod{3}$ . So either  $r^2$  or  $s^2$  are not  $0 \pmod{3}$ , and  $t = r^2 + s^2 \not\equiv 0 \pmod{3}$ . Thus,  $t$  and 3 are coprimes in  $\mathbb{Z}$  so there exists some  $u, v \in \mathbb{Z}$  such that  $3u + tv = 1$ .

Finally, note that  $r^2 + s^2 = (r + si)(r - si)$  so

$$\underbrace{3}_{\in I} u + \underbrace{(r + si)}_{\in I} (r - si)v = 1,$$

and it follows that  $1 \in I$ , and  $I = \mathbb{Z}[i]$ . Hence  $\langle 3 \rangle$  is indeed maximal, and it is also irreducible.

## LECTURE 6

### FIELDS

---

In the previous lecture we've discussed some new properties: primality and irreducibility. First we looked at prime elements and irreducible elements and afterwards we discussed prime ideals and maximal ideals. We discussed the strong connections between these definitions. Also, we have seen how we can construct new ideals using existing ideals and for dessert we started testing these properties in the Gaussian Integers ring.

## 6.1 Ideal Properties - Continued

### 6.1.1 From Previous Lecture

**Claim 1.** *let  $r \in R$  commutative ring*

1.  $r$  is prime  $\Leftrightarrow \langle r \rangle$  is prime
2.  $r$  is irreducible  $\Leftrightarrow \langle r \rangle$  is maximal ideal among principal ideals
3.  $\langle r \rangle$  is maximal  $\Rightarrow \langle r \rangle$  is prime
4.  $r$  is prime  $\Rightarrow r$  is irreducible
5.  $r$  is irreducible  $\Rightarrow \langle r \rangle$  is maximal if  $R$  is P.I.D
6.  $\langle r \rangle$  is maximal  $\Leftrightarrow \frac{R}{\langle r \rangle}$  is a field

### 6.1.2 Gaussian Integers

In the previous lecture, we saw that  $\langle 3 \rangle$  is maximal in  $\mathbb{Z}[i]$ , and that 5 is not irreducible since  $5 = (2 + i)(2 - i)$  and hence  $\langle 5 \rangle$  not maximal from (2).

**Claim 2.**  $\langle 2 + i \rangle$  is maximal in  $\mathbb{Z}[i]$

*Proof.* We could have shown like we did in  $\langle 3 \rangle$  that there is no bigger ideal, but let's try another way. Now we will show that  $\mathbb{Z}[i]/\langle 2 + i \rangle$  is a field and by (6) it follows that  $\langle 2 + i \rangle$  is maximal.

$$\langle 2 + i \rangle = \{(2 + i)(a + bi) \mid a, b \in \mathbb{Z}\} = \{(2a - b) + (a + 2b)i \mid a, b \in \mathbb{Z}\}$$

Let's define  $\alpha = 2a - b$ ,  $\beta = a + 2b$ . Let's pay attention:  $2\alpha + \beta = 2(2a - b) + (a + 2b) = 5a$

We can quickly see that a necessary property (but might not be sufficient) for elements in our ideal is that  $5 \mid (2\alpha + \beta)$ .



**Lemma 6.3.**  $\langle 2+i \rangle = \{a+bi \mid a, b \in \mathbb{Z}, 2a+b \equiv_5 0\}$  (Our property is also sufficient)

*Proof.* We only need to show that if  $2a+b \equiv_5 0$ ,  $a+bi \in \langle 2+i \rangle$

$$a+bi \in \langle 2+i \rangle \iff a+bi+ai(2+i) \in \langle 2+i \rangle \iff (2a+b)i \in \langle 2+i \rangle$$

We know that  $5 = (2+i)(2-i)$  so  $5 \in \langle 2+i \rangle$  so  $(2a+b)i \in \langle 2+i \rangle$  □

**Lemma 6.4.** Let  $I = \langle 2+i \rangle$ , we'll prove that  $\mathbb{Z}[i]/I = \{I, i+I, 2i+I, 3i+I, 4i+I\}$

*Proof.* First of all, we will show that every two elements are different. Assume  $ai+I = bi+I$  so  $ai-bi \in I$  so  $(a-b) \equiv_5 0$  so  $a=b$  (both are in  $0, 1, 2, 3, 4$ ).

Now we will show that there are no others. Let there be  $a+bi+I$  so  $a+bi+I = a+bi+ai(2+i)+I = a+bi+2ai-a = (b+2a)i+I$ .

Since  $5 \in I$ , we can choose whichever  $c$  we want such that  $b+2a-5c \in \{0, 1, 2, 3, 4\}$  and  $a+bi+I = (b+2a-5c)i+I$  and we get that:  $a+bi+I \in \{I, i+I, 2i+I, 3i+I, 4i+I\}$  □

**Lemma 6.5.**  $\mathbb{Z}[i]/I$  is a field

*Proof.* Since it is a quotient ring, it satisfies all the axioms of a field, except for the existence of multiplicative inverses. We will show that easily.

First of all  $2i+I$  is the unit element:

$$\text{Let } c \in \{0, 1, 2, 3, 4\} \text{ so } (2i+I)(ci+I) = -2c+I = -2c+c(2+i)+I = ci+I$$

Now we need to show that every element (except of zero which is  $I$ ) has an element such that their multiply result is  $2i+I$ .

$$(3i+I)(3i+I) = -9+I = -9+5+2(2+i)+I = 2i+I \text{ so } 3i+I \text{ is its own inverse.}$$

$$(4i+I)(i+I) = -4+I = -4+2(2+i)+I = 2i+I \text{ so } 4i+I \text{ is the inverse of } i+I.$$

So we showed that every element has an inverse. □

We proved that  $\mathbb{Z}[i]/I$  is a field. So it follows that  $I$  is maximal. □

### 6.1.3 $\mathbb{Z}[\sqrt{-5}]$

We can define a different ring like this:

$$\mathbb{Z}[\sqrt{-5}] = \{a+b\sqrt{-5} \mid a, b \in \mathbb{Z}\}$$

In this ring, some strange things happen. Let's look at two different factorizations of 6:  $6 = 2 * 3 = (1 + \sqrt{-5})(1 - \sqrt{-5})$

We'll recall that  $p$  is prime if  $p \mid ab \rightarrow p \mid a$  or  $p \mid b$ , but in this case we will see that  $2 \mid 6$  but  $2 \nmid (1 + \sqrt{-5})$  and  $2 \nmid (1 - \sqrt{-5})$ .

**Lemma 6.6.**  $2 \nmid (1 + \sqrt{-5})$  and  $2 \nmid (1 - \sqrt{-5})$

*Proof.* Assume by contradiction that  $2|(1 + \sqrt{-5})$ . So, there exists  $a, b \in \mathbb{Z}$  s.t.  $2(a + b\sqrt{-5}) = (1 + \sqrt{-5})$ .

It follows that,  $2a = 1$ , but  $a = \frac{1}{2} \notin \mathbb{Z}$ . The same is true for  $2 \nmid (1 - \sqrt{-5})$ .  $\square$

**Lemma 6.7.** *2 is irreducible*

*Proof.* We see that every element in our ring is also an element in the complex numbers field. Recall the norm of complex number:

$$\|a + b\sqrt{-5}\|_2 = \|a + b\sqrt{5}i\|_2 = \sqrt{a^2 + 5b^2}$$

Notice that every element  $0 \neq r \in \mathbb{Z}[\sqrt{-5}]$  in our ring it holds  $\|r\|_2 \geq 1$ . We are looking for two elements  $x, y \in \mathbb{Z}[\sqrt{-5}]$  which their multiplication result is 2. So, because  $\|x\|_2, \|y\|_2 \geq 1$ , it follows that  $\|x\|_2, \|y\|_2 \leq 2$ . If  $b \neq 0$ , its norm is at least  $\sqrt{5}$  so  $b_y = b_x = 0$ . Therefore we are in  $\mathbb{Z}$  and there is only  $2 = 1 * 2 = (-1) * (-2)$  where 1, -1 are units also in  $\mathbb{Z}[\sqrt{-5}]$   $\square$

Because 2 is irreducible but not prime it follows that  $\mathbb{Z}[\sqrt{-5}]$  isn't P.I.D. This factorization reminds us of a fact from  $\mathbb{Z}$ :  $60 = 15 * 4 = 10 * 6$  where  $4 \nmid 10, 6$  but we know that it is because we didn't factorize 60 to its prime components.

Recap some definitions:

**Definition 6.8.** Let  $R$  be commutative ring and  $I, J$  ideals in  $R$ .

$$\langle i, j \rangle = \langle i \rangle + \langle j \rangle = \{xi + yj | x, y \in R\}$$

$$IJ = \left\{ \sum_{k=1}^n i_k j_k \mid \forall n \in \mathbb{N}, i_k \in I, j_k \in J \right\}$$

Fact:  $IJ \subseteq I \cap J$

Let's take a look at  $\langle 6 \rangle$  and see the factorization ideals behavior. Obviously,  $\langle 6 \rangle = \langle 3 \rangle \langle 2 \rangle$

**Lemma 6.9.**  $\langle 3 \rangle = \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle$

*Proof.* Let's start by proving  $\langle 3 \rangle \subseteq \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle$ .

For that we want to show that the element 3 can be produced by a multiplication of  $ij$  where  $i \in \langle 3, 1 - \sqrt{-5} \rangle, j \in \langle 3, 1 + \sqrt{-5} \rangle$ . We'll notice from the fact above  $3 \in \langle 3, 1 - \sqrt{-5} \rangle \cap \langle 3, 1 + \sqrt{-5} \rangle$  doesn't necessarily mean that  $3 \in \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle$ . We know that  $9 = 3 * 3 \in IJ$ , and  $-6 = -(1 - \sqrt{-5})(1 + \sqrt{-5}) \in IJ$ , so from additive closure it follows that  $3 = 9 + (-6) \in IJ$

Let's now prove the other side:  $\langle 3 \rangle \supseteq \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle$

We'll show that  $\forall ij \in \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle, 3|ij$

let  $ij \in \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle. ij = (a * 3 + b * (1 - \sqrt{-5}))(c * 3 + d * (1 + \sqrt{-5})) =$

$$9ac + 3(1 + \sqrt{-5})ad + 3(1 - \sqrt{-5})bc + \overbrace{(1 - \sqrt{-5})(1 + \sqrt{-5})}^6 cd.$$

So, it's easy to see that  $3|9ac, 3|3(1 + \sqrt{-5})ad, 3|3(1 - \sqrt{-5})bc, 3|6cd$ , so  $3|ij$ .  $\square$

**Lemma 6.10.**  $\langle 2 \rangle = \langle 2, 1 + \sqrt{-5} \rangle^2$

*Proof.* Let's start by proving  $\langle 2 \rangle \subseteq \langle 2, 1 + \sqrt{-5} \rangle^2$

We want to show that the element 2 can be produced by a multiplication of  $ij$   $i, j \in \langle 2, 1 + \sqrt{-5} \rangle$

$-4 = -2 * 2 \in \langle 2, 1 + \sqrt{-5} \rangle^2$ ,  $6 = \overbrace{(2 - (1 + \sqrt{-5}))}^{\in \langle 2, 1 + \sqrt{-5} \rangle^2} \overbrace{(1 + \sqrt{-5})}^{\in \langle 2, 1 + \sqrt{-5} \rangle^2} \in \langle 2, 1 + \sqrt{-5} \rangle^2$  so from additive closure it follows that  $2 = 6 + (-4) \in \langle 2, 1 + \sqrt{-5} \rangle^2$

Let's prove the other side:  $\langle 2 \rangle \supseteq \langle 2, 1 + \sqrt{-5} \rangle^2$

We'll show that  $\forall ij \in \langle 2, 1 + \sqrt{-5} \rangle^2, 2|ij$

let  $ij \in \langle 2, 1 + \sqrt{-5} \rangle^2. ij = (a * 2 + b * (1 + \sqrt{-5})) \overbrace{(c * 2 + d * (1 + \sqrt{-5}))}^{-4 + 2\sqrt{-5}} = 4ac +$

$2(1 + \sqrt{-5})ad + 2(1 + \sqrt{-5})bc + (1 + \sqrt{-5})(1 + \sqrt{-5})cd.$

So, it's easy to see that  $2|4ac, 2|2(1 + \sqrt{-5})ad, 2|2(1 + \sqrt{-5})bc, 2|-2(2 + \sqrt{-5})cd,$  so  $2|ij.$  □

From the previous lemmas we get that

$$\langle 6 \rangle = \langle 3 \rangle \langle 2 \rangle = \langle 3, 1 - \sqrt{-5} \rangle \langle 3, 1 + \sqrt{-5} \rangle \langle 2, 1 + \sqrt{-5} \rangle \langle 2, 1 + \sqrt{-5} \rangle$$

Similarly, we can show that

$$\langle 3, 1 + \sqrt{-5} \rangle \langle 2, 1 + \sqrt{-5} \rangle = \langle 1 + \sqrt{-5} \rangle \text{ and}$$

$$\langle 3, 1 - \sqrt{-5} \rangle \langle 2, 1 + \sqrt{-5} \rangle = \langle 1 - \sqrt{-5} \rangle$$

so we also get  $\langle 6 \rangle = \langle 1 + \sqrt{-5} \rangle \langle 1 - \sqrt{-5} \rangle$

and because we saw earlier that  $6 = (1 + \sqrt{-5})(1 - \sqrt{-5})$ , we can see that the ideals behave correctly.

Earlier, we didn't see the connection between  $6 = 3 * 2 = (1 + \sqrt{-5})(1 - \sqrt{-5})$ , but when we worked with the ideals, they factorized "better" than the numbers.

## 6.2 Fields

### 6.2.1 Fields of Fractions

For every integral domain, we would like to find the "smallest" field that contains this domain. We have some intuition on how this field would look like (fractions) and when the two fractions are equal. Let's use this intuition to define it formally.

For integral domain  $R$  we'll define the set  $S = \{(a, b) \mid a \in R, b \in R \setminus \{0\}\}$

We'll define equivalence relation  $\sim$  on  $S$  s.t  $(a, b) \sim (c, d) \iff ad = bc$

$Q(R) = \{[(a, b)] \mid (a, b) \in S\}$  (  $[\ ]$  =equivalence class of  $\sim$  )

Lets define the  $(+, *)$  of  $Q(R)$ :

$$[(a, b)] * [(c, d)] = [(a * c, b * d)]$$

$$[(a, b)] + [(c, d)] = [(a * d + b * c, b * d)]$$

(For the multiplication we now use that this is integral domain and  $b * d \neq 0$ )

It's easy to see that this is well defined and that  $\forall r \in R$  there exists  $[(r, 1)] \in Q(R)$  which isomorphic to  $r$

In addition,  $\forall r \neq 0 \in Q(R)$  exists  $s \in Q(R)$  s.t  $r * s = [(1, 1)]$  (it is  $s = [(1, r)]$ ).

For example, if  $Q(\mathbb{Z}) = \mathbb{Q}$  and the equivalence relation  $\sim$  defines all the fractions we are familiar with  $[(a, b)] = \frac{a}{b} = \frac{c * a}{c * b}$

**Theorem 6.11.** *F is a field and  $R \subseteq F$ . So  $\exists \varphi : Q(R) \rightarrow F$  s.t  $\varphi$  is monomorphism. (without proof)*

This is a very strong theorem which claims that every field that contains our ring actually contains the field of the fraction we just defined. That means that this field is actually the "smallest".

## 6.2.2 The Polynomial Ring

Sometimes we know an element which isn't in our field, and we will find a way to "add it" by an equation. For example,  $i \notin \mathbb{R}$  but we know that it is the element which solves  $x^2 + 1 = 0$ . Sometimes we won't be able to add the element (for example  $\pi \notin \mathbb{Q}$  which is transcendental).

**Definition 6.12.**  $R[x] = \{\sum_{i=0}^n r_i x^i \mid r_i \in R, n \geq 0\}$

Where  $R[x]$  located in the following order?

Field < P.I.D < I.D < C.R

Notice that  $R[x]$  can't be smaller than  $R$  in the order.

1. If  $R$  is c.r,  $R[x]$  is c.r.
2. If  $R$  is I.D,  $R[x]$  is also I.D:

Let  $(ax^n + \dots), (bx^m + \dots)$  polynomials, since  $R$  is I.D,  $a, b \in R$  it follows that  $a * b \neq 0$  so  $(ax^n + \dots) * (bx^m + \dots)$  is not the zero polynomial.

3. If  $R$  is P.I.D,  $R[x]$  is I.D. Of course that  $R[x]$  is at least I.D, because if  $R$  is P.I.D, it also I.D.

Why isn't  $R[x]$  a P.I.D? For example,  $\mathbb{Z}$  is P.I.D but the ideal  $\langle 2, x \rangle$  in  $\mathbb{Z}[x]$  is not principal so  $\mathbb{Z}[x]$  is not P.I.D.

What if  $R$  is a field?

**Lemma 6.13.** *If  $R$  is a field,  $R[x]$  is P.I.D*

*Proof.* 1.  $R[x]$  is not a field

Because there is no inverse to  $x$  in  $R[x]$ .

2.  $R[x]$  is P.I.D

For each ideal, we'll find a polynomial that creates it.

Reminder: degree of polynomial -

- $\deg(a) = 0$
- $\deg(0) = -\infty$
- Let  $f, g \in R[x]$ ,  $\deg(fg) = \deg(f) + \deg(g)$
- Let  $f, g \in R[x]$ ,  $\deg(f + g) \leq \max(\deg(f), \deg(g))$

In order to prove that  $R[x]$  is P.I.D, we'll show that every  $\{0\} \neq I \subseteq R[x]$  ideal is principal.

We know that every polynomial  $g(x)$  can be written as  $g(x) = q(x)f(x) + r(x)$  where  $r(x)$  is a remainder such that  $\deg(r(x)) < \deg(f(x))$  or  $r(x) = 0$

We'll take  $f(x) \neq 0$  with minimal degree in  $I$ .

$\underbrace{g(x)}_{\in I} = q(x)\underbrace{f(x)}_{\in I} + r(x)$ , so it follows that  $q(x)f(x) \in I$  and that  $r(x) \in I$ . So  $r(x) = 0$ , because  $f(x)$  has minimal degree in  $I$ .

Therefore,  $I = \langle f(x) \rangle$  □

**Definition 6.14.** A polynomial is monic if its leading coefficient is 1. ( $a_n = 1$ )

Also if  $I = \langle f(x) \rangle$  for  $f(x)$  monic,  $f(x)$  is unique.

## 6.2.3 Building Finite Fields

Let's look at an interesting way to build finite fields.

We just proved that if  $F$  is a field,  $F[x]$  is P.I.D We also remember that in a P.I.D for every irreducible element  $f$ ,  $\langle f \rangle$  is a maximal ideal. The last crucial claim we proved is that for  $I$  maximal ideal,  $F[x]/\langle I \rangle$  is a field.

By looking at the quotient ring of infinite field and a maximal ideal, we can build a finite field. Let's look at some examples.

### 6.2.3.1 Finite field of 4

$F_2[x]$  is an infinite field (polynomials with coefficient 0 or 1).

We'll prove that  $F_2[x]/\langle x^2 + x + 1 \rangle$  is a field by showing that  $x^2 + x + 1$  is irreducible.

Let's try to factorize  $x^2 + x + 1$  into two polynomials which are not the units.

Let  $q(x), f(x)$  polynomials such that  $x^2 + x + 1 = q(x)f(x)$ . It follows that,  $\deg(q(x)) + \deg(f(x)) = 2$ . The only polynomials with degree 1 are  $x, x + 1$ .

But,  $x(x+1) = x^2 + x$ ,  $x * x = x^2$ ,  $(x+1)(x+1) = x^2 + 1$  and they are all not equal to  $x^2 + x + 1$ .

Therefore  $x^2 + x + 1$  is irreducible, and  $\langle x^2 + x + 1 \rangle$  is maximal ( $F_2[x]$  is P.I.D), so it follows that  $F_2[x]/\langle x^2 + x + 1 \rangle$  is a field.

We'll show that  $|F_2[x]/\langle x^2 + x + 1 \rangle| = 4$

**Lemma 6.15.**  $F_2[x]/\underbrace{\langle x^2 + x + 1 \rangle}_I = \{I, 1 + I, x + I, x + 1 + I\}$

We'll give an intuitive explanation why there are no more elements in this field:

Since  $I$  is the zero element in  $F_2[x]/I$ , we can think of  $x^2 + x + 1 = 0$ . From this we see that  $x^2 = x + 1$  ( $-1 = 1$  in  $F_2$ ).

So, for every degree  $d \geq 2$   $x^d$  can be written as a linear expression of  $x$ .

For example,  $x^3 = x^2 * x = (x + 1)x = x^2 + x = x + 1 + x = 1$ .

Now, we'll show that each element in this field is different:

We know that  $f + I = g + I \iff f - g \in I$

Assume by contradiction that two of the co-sets are equal,  $g + I = f + I$  such that  $\deg(g), \deg(f) \leq 1$  so by the previous sentence  $g - f \in I$  and we know that  $\deg(g - f) \leq 1$  but it means that  $x^2 + x + 1 | g - f$  but polynomial of degree 1 can't be factorized into polynomial of degree 2.

Define  $I = 0$ ,  $1 + I = 1$ ,  $x + I = a$ ,  $x + 1 + I = b$

Let's explore the multiplication matrix of this field:

	0	1	a	b	
0	0	0	0	0	$a * a = (x + I)(x + I) = x^2 + I = x + 1 + I = b$
1	0	1	a	b	$a * b = (x + I)(x + 1 + I) = x^2 + x + I = 1 + I = 1$
a	0	a	b	1	$b * b = (x + 1 + I)(x + 1 + I) = x^2 + 1 + I = x + I = a$
b	0	b	1	a	

### 6.2.3.2 Finite field of 8

We'll take  $F_2[x]/\langle x^3 + x + 1 \rangle$ . One can show that  $x^3 + x + 1$  is irreducible (similar to the previous example), Therefore  $F_2[x]/\langle x^3 + x + 1 \rangle$  is a field.

$F_2[x]/\underbrace{\langle x^3 + x + 1 \rangle}_I = \{I, 1 + I, x + I, x + 1 + I, x^2 + I, x^2 + 1 + I, x^2 + x + I, x^2 + x + 1 + I\}$

Same as in the previous example, these are the only elements in the field because we can't factorize a polynomial with a polynomial with a higher degree.

### 6.2.3.3 Finite field of 9

The same procedure can be done with  $F_3[x]/\langle x^2 + 1 \rangle$  and we will get a finite field of size 9.

In general, we can build a finite field of size  $p^n$  for  $p$  prime, by taking  $F_p[x]$  and an irreducible polynomial  $f(x)$ ,  $\deg(f(x)) = n$ . The field will be  $F_p[x]/\langle f(x) \rangle$

## 6.2.4 Building The Complex Field

This might be the most exciting part of this subject - building  $\mathbb{C}$ . In  $\mathbb{R}$  the polynomial  $x^2 + 1$  has no roots and we know it is irreducible. Therefore, we can look at  $\mathbb{R}[x]/\langle x^2 + 1 \rangle$  when we know it is a field.

$$\mathbb{R}[x]/\langle x^2 + 1 \rangle = \{f(x) + \langle x^2 + 1 \rangle \mid f(x) \in \mathbb{R}[x]\} = \{a + bx + \langle x^2 + 1 \rangle \mid a, b \in \mathbb{R}\}$$

Similar to the previous examples, we can look only at  $a + bx$  where we know " $x^2 + 1 = 0$ " or equally " $x^2 = -1$ ".

We'll look at our polynomial in the new field, and call it  $t^2 + 1$ . We can actually see it has a root - the co-set:  $x + I$  (where  $I = \langle x^2 + 1 \rangle$ ). We use  $t$  since  $x$  is no longer an abstract symbol but an element in our field.

We'll substitute  $t = x + I$  in the polynomial  $t^2 + 1$ :

$$(x + I)^2 + 1 = x^2 + I + 1 = x^2 + 1 + I = I \text{ since } (x^2 + 1) \in I$$

Let's explore how multiplication works here:

$$(ax + b + I)(cx + d + I) = acx^2 + (bc + ad)x + bd + I = (bc + ad)x + (bd - ac) + I$$

This is actually similar to multiplying elements in  $\mathbb{C}$ ...

# LECTURE 7

## FIELD EXTENSIONS

---

### 7.1 Recap - Constructing $\mathbb{C}$ from $\mathbb{R}$

We want to find roots for the polynomial  $p(x) = x^2 + 1$  in some other related field. We look at the following quotient ring:

$$\overline{\mathbb{C}} = \mathbb{R}[x] / \underbrace{\langle x^2 + 1 \rangle}_I$$

Note that  $\mathbb{R}[x]$  is a ring and that  $x^2 + 1$  is an irreducible element in that ring. Since  $\mathbb{R}[x]$  is a p.i.d,  $\langle x^2 + 1 \rangle$  is a maximal ideal, thus  $\overline{\mathbb{C}}$  is a field.

**Claim 1.**  $\mathbb{R} \hookrightarrow \overline{\mathbb{C}}$ , i.e  $\mathbb{R}$  is embedded in  $\overline{\mathbb{C}}$ , that is -

$$\exists \varphi : \mathbb{R} \rightarrow \overline{\mathbb{C}}$$

s.t  $\varphi$  is a monomorphism (1:1)

*Proof.* We choose  $\varphi$  to be  $\varphi(r) = r + I$ .

- $\varphi$  is a homomorphism:
  - $\varphi(1_{\mathbb{R}}) = 1 + I = 1_{\overline{\mathbb{C}}}$
  - $\varphi(r + s) = (r + s) + I = (r + I) + (s + I) = \varphi(r) + \varphi(s)$
  - $\varphi(r \cdot s) = (r \cdot s) + I = (r + I) \cdot (s + I) = \varphi(r) \cdot \varphi(s)$
- $\varphi$  is a monomorphism since  $\ker \varphi = \{0\}$ .

□

**Claim 2.**  $y^2 + 1$  has a root in  $\overline{\mathbb{C}}$



Note that even though  $y^2 + 1$  and  $x^2 + 1$  are the same polynomial (as they have the exact same coefficients),  $y^2 + 1$  lies in  $\overline{\mathbb{C}}[y]$ , and should be written as following:

$$\begin{aligned} \overline{\mathbb{C}}[y] \ni \quad & \varphi(1)y^2 + \varphi(1) = \\ & (1 + I)y^2 + (1 + I) = \\ & (1 + \langle x^2 + 1 \rangle)y^2 + (1 + \langle x^2 + 1 \rangle) \end{aligned}$$

*Proof.* We will show that  $x + I$  is the root of  $y^2 + 1 \in \overline{\mathbb{C}}[x]$ .  
Let's assign  $x + I$  to  $y$  in the above polynomial:

$$\begin{aligned} & \underbrace{(1 + I)}_{\substack{\text{An element} \\ \text{of the} \\ \text{quotient ring}}} \cdot \underbrace{(x + I)^2}_{\substack{\text{An element} \\ \text{of the} \\ \text{quotient ring,} \\ \text{squared}}} + \underbrace{1 + I}_{\substack{\text{An element} \\ \text{of the} \\ \text{quotient ring}}} = \\ & 1 \cdot x \cdot x + I \quad + 1 + I = \\ & x^2 + 1 + I = \\ & 0 + I \end{aligned}$$

□

**Corollary 7.3.**  $i$  is the root of  $x^2 + 1$  in  $\mathbb{C}$ .

We saw that  $x$  is the root of  $y^2 + 1$ , and we can change our markings.

## 7.2 Field Extensions

**Definition 7.4.** Let  $F, K$  be fields.  $F \subset K$  is a *field extension* w.r.t  $K$ 's addition and multiplication operations -  $(+, \cdot)$ .

That is,  $F \subset K^{(+, \cdot)}$  is a field extension if by using  $K$ 's addition and multiplication operations we create the field  $(F, +, \cdot)$

If  $F \subset K$  is a field extension we say:

- $F$  is a *subfield* of  $K$
- $K$  is an *extension* (or a *field extension*) of  $F$

We mark a field extension with  $K/F$

**Examples:**

- $\mathbb{R}$  is an extension of  $\mathbb{Q}$
- $\mathbb{C}$  is **not** an extension of  $\mathbb{Z}_2$ , addition in  $\mathbb{C}$  is not like addition in  $\mathbb{Z}_2$ .
- $\mathbb{C}$  is an extension of  $\mathbb{R}$ . Recall that  $\overline{\mathbb{C}} = \mathbb{R}[x] / \langle x^2 + 1 \rangle$ , this means addition in  $\mathbb{C}$  has the following form:

$$(f(x) + \langle x^2 + 1 \rangle) + (g(x) + \langle x^2 + 1 \rangle)$$

We assign  $r_1, r_2 \in \mathbb{R}$ :

$$\begin{aligned} (r_1 + \langle x^2 + 1 \rangle) + (r_2 + \langle x^2 + 1 \rangle) = \\ r_1 + r_2 + \langle x^2 + 1 \rangle \end{aligned}$$

The result in  $\mathbb{R}$  is the representative  $r_1 + r_2$ .

**Observation.** If  $K/F$  is a field extension then, in particular,  $K$  is a *vector space* over  $F$ .

To construct this vector space we create linear combinations of elements from  $K$  with coefficients from  $F$  (we don't multiply elements from  $K$  with each other), for example:

$$\begin{aligned} k_1 + k_2 \\ f k \\ f_0 f_1 k_1 + f_2 k_2 \end{aligned}$$

**Definition 7.5.** Let  $F, K$  be fields s.t.  $K/F$  is a field extension. we define the *degree extension* of  $K/F$  by  $[K : F] := \dim_F K$ , which is equal to the minimum size of a group  $\{k_1, \dots, k_n\} \subseteq K$  s.t.  $\forall k \in K \exists f_1, \dots, f_n \in F : k = \sum_{i=0}^n f_i k_i$

**Claim 6.**  $[\mathbb{C} : \mathbb{R}] = 2$

*Proof.* Using  $1, i \in \mathbb{C}$  we can represent all members in  $\mathbb{C}$  as a linear combination with prefixes from  $\mathbb{R}$ :

$$\forall c \in \mathbb{C} \exists a, b \in \mathbb{R} : c = a + ib \rightarrow [\mathbb{C} : \mathbb{R}] \leq 2,$$

in addition they are linearly independent - if there exist  $0 \neq a, b \in \mathbb{R}$  s.t.  
 $a \cdot 1 + i \cdot b = 0$  we will get that  $i = -\frac{a}{b} \in \mathbb{R}$ , which contradicts the fact that  
 $i \notin \mathbb{R} \rightarrow [\mathbb{C} : \mathbb{R}] = 2$  □

**Claim 7.**  $[\mathbb{R} : \mathbb{Q}] = \infty$

**Claim 8.**  $[\mathbb{F}_4 : \mathbb{F}_2] = 2$

*Proof.* By definition,

$$\mathbb{F}_4 = \mathbb{F}_2[x] / \langle x^2 + x + 1 \rangle = \{ax + b + \langle x^2 + x + 1 \rangle \mid a, b \in \mathbb{Z}_2\}$$

using  $1, x \in \mathbb{F}_4$  we can represent all the members in  $\mathbb{F}_4$ :

$$\forall f \in \mathbb{F}_4 \exists a, b \in \mathbb{F}_2 : f = a \cdot 1 + b \cdot x \rightarrow [\mathbb{F}_4 : \mathbb{F}_2] \leq 2$$

and since  $1, x$  are linearly independent,  $[\mathbb{F}_4 : \mathbb{F}_2] = 2$ . □

**Claim 9.**  $[\mathbb{F}_8 : \mathbb{F}_2] = 3$

*Proof.* By definition,

$$\mathbb{F}_8 = \mathbb{F}_2[x] / \langle x^3 + x + 1 \rangle = \{ax^2 + bx + c + \langle x^3 + x + 1 \rangle \mid a, b, c \in \mathbb{Z}_2\}.$$

using  $1, x, x^2 \in \mathbb{F}_8$  we can represent all the members in  $\mathbb{F}_8$ :

$$\forall f \in \mathbb{F}_8 \exists a, b, c \in \mathbb{F}_2 : f = a \cdot x^2 + b \cdot x + c$$

therefore  $[\mathbb{F}_8 : \mathbb{F}_2] \leq 3$ , and since  $1, x, x^2$  are linearly independent  $[\mathbb{F}_8 : \mathbb{F}_2] = 3$ . □

**Claim 10.**  $[\mathbb{F}_8 : \mathbb{F}_4]$  is undefined since  $\mathbb{F}_8$  is not a field extension of  $\mathbb{F}_4$ :  $x \cdot_{\mathbb{F}_8} x = x^2 \notin \mathbb{F}_4$ .

**Definition 7.11.** Let  $F$  be a field. the characteristic (char) of  $F$  is the least integer  $n \geq 1$  s.t.  $\underbrace{1 + \dots + 1}_{n \text{ times}} = 0$  if such  $n$  exists, otherwise the char is defined to be 0.

**Claim 12.** If  $F$  is a field of char 0, then  $\mathbb{Q} \hookrightarrow F$ .

*Proof.*  $F$  is a field therefore  $\forall a \in F$  exists both  $-a$  and  $a^{-1}$  s.t.  $a + (-a) = 0_F$  and  $a \cdot a^{-1} = 1_F$  we will build a monomorphism  $\varphi: \mathbb{Q} \rightarrow F$  in the following way:

$$\varphi\left(\frac{a}{b}\right) = ab^{-1}$$

- $\varphi$  is homomorphism since  $\forall \frac{a}{b}, \frac{c}{d} \in \mathbb{Q}$ :
  - $\varphi(1_{\mathbb{Q}}) = \varphi(\frac{m}{m}) = m_F m_F^{-1} = 1_F \quad (m \neq 0)$
  - $\varphi(\frac{a}{b} \cdot \frac{c}{d}) = \varphi(\frac{ac}{bd}) = a_F c_F d_F^{-1} b_F^{-1} = (a_F b_F^{-1}) \cdot (c_F d_F^{-1}) = \varphi(\frac{a}{b}) \cdot \varphi(\frac{c}{d})$
  - $\varphi(\frac{a}{b} + \frac{c}{d}) = \varphi(\frac{ad+bc}{bd}) = (a_F d_F + b_F c_F)(b_F d_F)^{-1} = a_F d_F d_F^{-1} b_F^{-1} + b_F c_F d_F^{-1} b_F^{-1} = a_F b_F^{-1} + c_F d_F^{-1} = \varphi(\frac{a}{b}) + \varphi(\frac{c}{d})$
- $\varphi$  is a monomorphism since:
 
$$\varphi(\frac{a}{b}) = 0 \iff a_F b_F^{-1} = 0 \iff a \cdot 1_F = 0 \iff a = 0 \rightarrow \ker \varphi = \{0\}.$$

□

**Claim 13.** *If  $F$  is a field of char  $n \neq 0$ , then  $n$  is prime.*

*Proof.* Assume  $n$  is not prime, therefore  $\exists a, b \in \mathbb{N}$  s.t.  $1 < a, b < n$  and  $n = ab$ . by the definition of fields characteristic,  $n \geq 1$  is the minimal integer s.t.  $n \cdot 1_F = 0$  and therefore  $a_F, b_F \neq 0_F$ . since  $a \cdot b = n$ ,  $a_F \cdot b_F = 0_F$ ,  $a_F$  and  $b_F$  are both zero divisors in  $F$ , but fields does not contain zero divisors since they are also integral domains, contradiction.

□

**Claim 14.** *If  $F$  is a field of prime char  $p \neq 0$ , then  $\mathbb{Z}_p \hookrightarrow F$ .*

*Proof.* We will build a monomorphism  $\varphi : \mathbb{Z}_p \rightarrow F$ .

$$\varphi(a) = a_F$$

- $\varphi$  is a homomorphism since  $\forall a, b \in \mathbb{Z}_p$ :
  - $\varphi(1_{\mathbb{Z}_p}) = 1_F$
  - $\varphi(a + b) = (a + b)_F = a_F + b_F = \varphi(a) + \varphi(b)$
  - $\varphi(a \cdot b) = (ab)_F = a_F \cdot b_F = \varphi(a) \cdot \varphi(b)$
- $\varphi$  is a monomorphism since  $\ker \varphi = \{0\}$ , otherwise there exists  $0 \neq n \in \mathbb{Z}_p$  s.t.  $\varphi(n) = n_F = 0$ , which contradicts the minimality of  $\text{char}(F) = p$ .

□

**Observation.** For every field  $F$ ,  $\mathbb{Z}_p \hookrightarrow F$  if  $\text{char}(F) = p \neq 0$ , and  $\mathbb{Q} \hookrightarrow F$  otherwise.

**Claim 15.** *Let  $F$  be a field.  $\text{char}(F)$  is 0 or a prime.*

*Proof.* Let's assume  $\text{char}(F) = n$  which is not a prime.

$n$  is not a prime, by definition  $\exists_{a,b}$  s.t  $n = a \cdot b$  and  $1 < a < n, 1 < b < n$  ( $a, b \neq 0$  since  $n$  is minimal).

$$\text{char}(F) = \underbrace{(1 + 1 + \dots + 1)}_n = \underbrace{n}_{\substack{\text{definition} \\ \text{of} \\ \text{char}}} = 0$$

thus,  $0 = a \cdot b$  which means  $a, b$  are zero divisors in  $F$  which is a contradiction to the assumption that  $F$  is a field.  $\Rightarrow \text{char}(F)$  is a prime or 0 □

**Example 6** can't be characteristic of any field, otherwise

$$1 + 1 + 1 + 1 + 1 + 1 = 0 \Rightarrow \underbrace{(1 + 1)}_2 + \underbrace{(1 + 1)}_2 + \underbrace{(1 + 1)}_2 = 0 \Rightarrow 2 \cdot 3 = 0 - \text{zero divisors}$$

**Definition 7.16.** Let  $F$  be a field and  $\varphi : \mathbb{Z} \rightarrow F$  homomorphism where  $\varphi(1) = 1_F$  we know that  $\ker \varphi = n\mathbb{Z} = \langle n \rangle$  or  $\{0\}$  because  $\mathbb{Z}$  is a p.i.d. From the first isomorphism theorem:

$$\mathbb{Z} / \ker \varphi \cong \text{Im} \varphi$$

We note that  $\text{Im} \varphi \subseteq F$  so it must be at least Integral Domain - meaning  $\mathbb{Z} / \ker \varphi$  is also Integral Domain. we know that  $R/P$  is I.D iff  $R$  is a ring and  $p$  is a prime ideal. in our case  $\mathbb{Z}$  is a ring and  $\text{Im} \varphi \subseteq F$  is I.D so  $\ker \varphi$  is a prime ideal. meaning  $\ker \varphi = p\mathbb{Z}$ . We define:

$$\underline{\text{char}}_2(F) := \begin{cases} p & \ker \varphi = p\mathbb{Z} \\ 0 & \ker \varphi = \{0\} \end{cases}$$

**Claim 17.**  $\text{char}_2(F) = \text{char}(F)$

*Proof.* Let  $\text{char}_2(F) = p$ . By definition,  $p \in \ker \varphi$ , meaning  $\varphi(p) = 0$

$$\underbrace{(1 + 1 + \dots + 1)}_p = 0 \Rightarrow p = \text{char}(F). \quad \square$$

**Corollary 7.18.** Every finite field  $K$  is of size  $p^n$  where  $p$  is a prime and  $n \in \mathbb{N}, n \geq 1$ .

*Proof.*  $\text{char}(K) \neq 0$  because  $K$  is a finite field.

We know  $K$  is a field with prime characteristic  $p$  so it contains  $\mathbb{F}_p: \mathbb{F}_p \hookrightarrow K$ .

For a finite  $K$  where  $\text{char}(K) = p$  then  $K$  is a vector space over  $\mathbb{F}_p$  with finite dimension because  $K$  is finite.

$$\dim_{\mathbb{F}_p} K = n < \infty$$

Therefore,  $\exists_{v_1, v_2, \dots, v_n \in K}$  s.t.  $\forall_{k \in K} \exists_{\lambda_1, \lambda_2, \dots, \lambda_n \in \mathbb{F}_p}$  s.t.

$$k = \sum_{i=1}^n \lambda_i v_i$$

$\{v_1, \dots, v_n\}$  are the basis vector set of  $K$  over  $F \Rightarrow$  thus  $|K| = p^n$ . □

### Example

$$\mathbb{F}_4 = \mathbb{F}_2[x] / \langle x^2 + x + 1 \rangle = \{a + bx + \langle x^2 + x + 1 \rangle \mid a, b \in \mathbb{F}_2\}$$

$$1 + 1 = 0 \Rightarrow \text{char} \mathbb{F}_4 = 2$$

$$\{1 + I, x + I\} \text{ are the basis vectors } \Rightarrow \dim_{\mathbb{F}_2} \mathbb{F}_4 = 2$$

$$|\mathbb{F}_4| = 2^2 = 4$$

**Claim 19.** If  $L \setminus K$  and  $K \setminus F$  are field extensions then

$$[L : F] = [L : K][K : F]$$

*Proof.* Let's assume first that  $[L : K]$  and  $[K : F]$  are finite.

$$[L : K] = n \quad [K : F] = m$$

$$\dim_K L = n \Rightarrow \exists_{\{l_1, \dots, l_n\}} \text{ vector basis of } L \text{ over } K$$

$$\dim_F K = m \Rightarrow \exists_{\{k_1, \dots, k_m\}} \text{ vector basis of } K \text{ in } F$$

$$\forall_{l \in L} \exists_{k'_1, \dots, k'_n \in K} \rightarrow l = \sum_{i=1}^n k'_i \cdot l_i$$

$$\forall_{k'_i \in K} \exists_{f_1^i, \dots, f_m^i \in F} \rightarrow k'_i = \sum_{j=1}^m f_j^i \cdot k_j$$

$$\forall_{l \in L} \rightarrow l = \sum_{i=1}^n \sum_{j=1}^m f_j^i \cdot k_j \cdot l_i$$

$$\text{which means that } V = \{v_{i,j} \mid v_{i,j} = k_j \cdot l_i, 1 < j < m, 1 < i < n\}$$

is a spanning set of  $L$  over  $F$

$$\dim_F L \leq n \cdot m$$

now we need to prove that the spanning set is also linearly independent. let's assume

it's linearly dependent.

$$\begin{aligned} \exists_{v_{i',j'} \in V} &\rightarrow \exists_{f_1, \dots, f_{n \cdot m} \in F} \rightarrow \sum_{\substack{v \in V \\ v \neq v_{i',j'}}} f_i v_i = v_{i',j'} \\ k_{j'} \cdot l_{i'} &= \sum_{\substack{i=1 \\ i \neq i'}}^n \sum_{\substack{j=1 \\ j \neq j'}}^m f_j^i \cdot k_j \cdot l_i \\ l_{i'} &= \sum_{\substack{l \in L \\ l \neq l_{i'}}} \frac{f_j^i \cdot k_j}{k_{j'}} \cdot l_i \end{aligned}$$

which is a contradiction of the vector basis of  $L$  over  $K$ .

if  $[L : K]$  or  $[K : F]$  are  $\infty \rightarrow [L : F] = \infty$ , if the latter is finite we could find a finite vector basis to both  $K$  over  $F$  and  $L$  over  $K$  in the same way we built the basis vector of  $L$  over  $F$ .  $\square$

**Definition 7.20.** Let  $K/F$  be a field extension.  $a \in k$  is algebraic over  $F$  if there exists a polynomial  $0 \neq f(x) \in F[x]$  s.t.  $f(a) = 0$

If such polynomial doesn't exist we define  $a$  to be transcendental over  $F$ .

**Example**  $i$  is algebraic over  $\mathbb{R}$  because

$$\mathbb{R}[x] \ni f(x) = 1 \cdot x^2 + 1$$

$$f(i) = 1 \cdot i^2 + 1 = 0$$

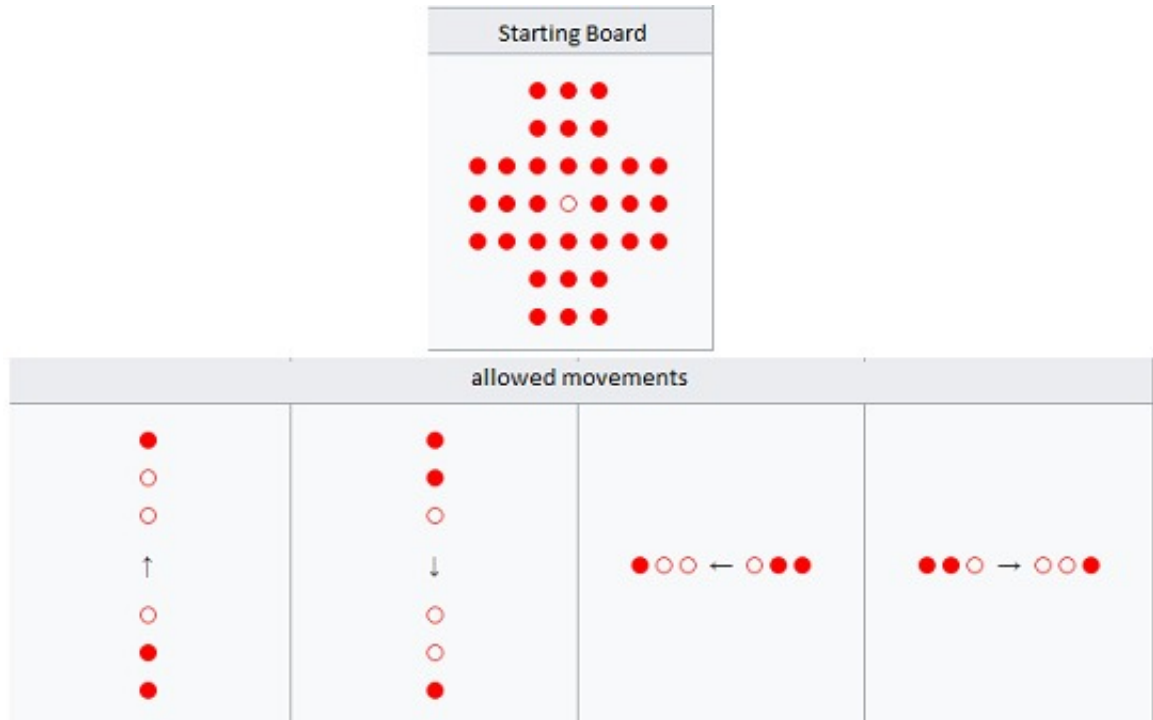
**Claim 21.** *The numbers  $\pi, e$  are transcendentals over  $\mathbb{Q}$ .*

**Note**  $\forall_{a \in F} \Rightarrow a$  is algebraic over  $F$ , because:

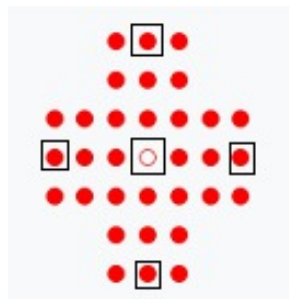
$$f(x) = x - a \quad f(x) \in F[x]$$

$$f(a) = a - a = 0$$

**Peg solitaire** Peg solitaire (or Solo Noble) is a board game for one player involving movement of pegs on a board with holes. The board starts filled with pegs beside the center, and using the allowed moves the player needs to remain with only one peg.



Using what we learned we can prove that the remaining peg can only be placed on one of the marked places





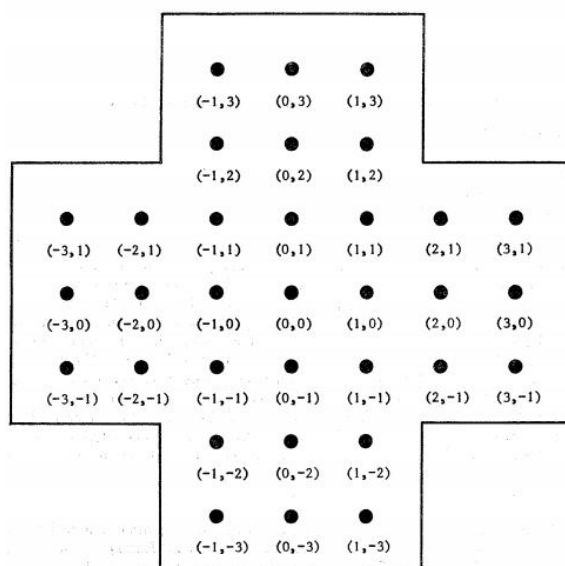
# LECTURE 8

## FIELDS AND POLYNOMIALS

---

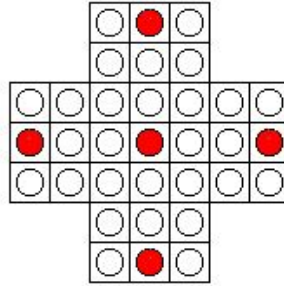
### 8.1 Proving final locations for a winning Peg Solitaire game

We finished the previous lecture starting to prove a property of Peg Solitaire. Peg solitaire (or Solo Noble) is a board game for one player involving movement of pegs on a board with holes. The game starts with the following board configuration of pegs. A dark dot at location  $(i,j)$  means that there is a peg at that coordinate.



**image 1:** starting configuration of the image of starting configuration.

A valid move is to jump a peg orthogonally over an adjacent peg into a hole two positions away and then to remove the jumped peg. The game is won whenever there is only one peg left on the board. It turns out, that the location of that winning peg can only be in one of 5 coordinates - the center of the board  $(0,0)$ , or the centers of the outer sides of the board  $(0,3)$ ,  $(3,0)$ ,  $(0,-3)$ ,  $(-3,0)$ , as shown in the image below.



**image 2:** possible locations for the last peg in a winning board.

We will now prove that the location of the final soldier in a winning board is always one of the those locations specified above. The idea is to understand the connection between the possible movements of each peg and a finite field of size 4:  $F_4 = F_2[x] / \langle x^2 + x + 1 \rangle = \{0, 1, x, x + 1\}$ . For every board configuration B, we shall define two functions and show that they maintain a certain invariant, summing up all pegs that are currently on board at all locations (i,j), under  $F_4$ .

The first function:

$$f(B) = \sum_{(i,j)} x^{i+j}, \quad f(B) \in F_4$$

The second function:

$$g(B) = \sum_{(i,j)} x^{i-j}, \quad g(B) \in F_4$$

both are sums going over the locations (i,j) of pegs for a board configuration B.

**Claim 1.** For a starting game configuration  $B_0$ :

$$f(B_0) = 1 = g(B_0)$$

*Proof.* In order to sum the result over the entire board, we will calculate the sum of every three holes. Every row of three holes with three pegs in them can be described (using  $f(B)$ 's phrasing) as:

$$x^c \quad x^{c+1} \quad x^{c+2}$$

•   •   •

So for that triplet of holes, if there are three pegs, then their accumulation

$$x^c + x^{c+1} + x^{c+2}$$

is added to  $f(B_0)$  for every triplet of pegs. Note that:

$$x^c + x^{c+1} + x^{c+2} = x^c(1 + x^1 + x^2) = x^c \cdot 0_{F_4} = 0$$

This can also be proven to each column of three holes. In total, in a starting configuration the board is full of triplets of pegs, with only one triplet left that includes the center hole, with no peg. The addition of that triplet to  $f$  is 1, and so in total we get  $f(B_0) = 1$ . The same can be proven easily for  $g(B_0)$   $\square$

**Claim 2.** *For every board configuration  $B$ :*

$$f(B) = 1 = g(B)$$

*Proof.* Assume a move in which a peg jumps over another peg adjacent to it in the row, into an empty hole according to the rules. The score of the triplet of coordinates for function  $f$  before the jump:

$$x^c, x^{c+1}, 0$$

$$\bullet \quad \bullet \quad o$$

(the pegs in the first and second locations are given scores, and the third hole is empty so no score is given). The score of the triplet of coordinates after the jump:

$$0, 0, x^{c+2}$$

$$o \quad o \quad \bullet$$

(There is only one peg left, in the hole that was first empty).

From the previous claim we see

$$x^c + x^{c+1} + x^{c+2} = 0 \Rightarrow x^c + x^{c+1} = x^{c+2}$$

( $F_4$  is binary). This can be proven for the  $g$  invariant as well, and so every board configuration gets the same score both before and after the move.  $\square$

Using the two claims above, we conclude that in a winning configuration  $B_f$ , the sum

of both  $f, g$  remains 1. There is only one soldier left, in location  $(i, j)$ . This means:

$$f(B_f) = x^{i+j} = 1, \quad g(B_f) = x^{i-j} = 1$$

We are in  $F_4$ , and if we mark  $i+j=c$ ,  $i-j=d$ , we see that the only option for  $x^c = 1 = x^d$  is for 3 to divide both  $c$  and  $d$ :  $3|c, 3|d$ . This only happens for coordinates  $(0,3), (3,0), (0,-3), (-3,0)$  and  $(0,0)$ , which are always the final coordinates for a winning peg in Peg Solitaire.

## 8.2 Back to Fields

**Definition 8.3.** Let  $F \subseteq K$  and  $K$  is a field extension of  $F$ . We say that  $a \in K$  is **algebraic over  $F$**  if there exists  $f(x) \in F[X]$  such that  $f(x) \neq 0$  and  $f(a) = 0$ . If  $a$  isn't algebraic over  $F$ , then we say that  **$a$  is transcendental**.

There is a small issue with the definition above because we think of the polynomial as a function. In order to use that definition, we have to define assignment in a polynomial.

A polynomial is an infinite sequence of coefficients in  $F$  that has a finite number of coefficients that are not zero:  $(f_0, f_1, \dots) \in F$ . Each of the coefficients corresponds to  $x^0, x^1, \dots$  and so we get

$$f_0x^0 + f_1x^1 + \dots$$

When assigning  $a$ , we get:

$$f_0 \cdot a^0 + f_1 \cdot a + \dots + f_n \cdot a^n$$

But in doing so, we think of the polynomial as a function. We need to define this. The way to define assignment inside a polynomial is done using homomorphism (for example, assigning  $m$ ):

$$ev_m : R[m] \rightarrow R$$

i.e.

$$a_0 + a_1 \cdot x + \dots + a_n x^n \rightarrow a_0 \cdot m^0 + a_1 \cdot m^1 + \dots + a_n \cdot m^n$$

from this definition for polynomials, we can conclude:

1.  $f(m) + g(m) = (f + g)(m)$
2.  $f(m)g(m) = (f \cdot g)(m)$ .

Since it is a homomorphism, we know it has a kernel:

$$\text{Im } ev_m = R$$

and

$$\text{Ker } ev_m = \langle x - m \rangle$$

i.e. all the polynomials that  $m$  divides. (In other words, if we assign  $m$  and the result of the assignment is 0, then  $x - m$  divides the polynomial. Thus, from the first homomorphism theorem we get:

$$R[x] / \langle x - m \rangle \simeq R$$

Assigning  $m$  in a polynomial  $f$  is in fact, finding the representative with the lowest degree in the coset of the polynomial:  $f(x) + \langle x - m \rangle = r + \langle x - m \rangle$ .

**Definition 8.4.** Suppose  $K$  is a field extension over  $F$ .  $K$  is called an **algebraic extension** if for every  $a \in K$ ,  $a$  is algebraic over  $F$ .

**Definition 8.5.** Suppose  $K$  is a field extension over  $F$ .  $K$  is said to be a **finite extension** if

$$[K : F] < \infty$$

(Reminder:  $[K : F] = \dim_F K$ )

**Theorem 8.6.** *Every finite field extension is an algebraic extension.*

*Proof.*  $K$  is a field extension over  $F$  such that  $[K : F] = n < \infty$ . We need to show that the extension is algebraic. So, Let some  $a \in K$ . We want to show that there exists  $f(x) \in F[X]$  s.t.  $f(a) = 0$ .

Let us examine the following  $n+1$  elements:  $1, a, a^2, \dots, a^n \in K$ . Because we selected more than  $n$  elements while there are exactly  $n$  elements in the basis, we deduce that the elements are linearly dependent:  $\exists f_0, f_1, \dots, f_n \in F$  not all zeros s.t.

$$f_0 \cdot 1 + f_1 \cdot a + \dots + f_n \cdot a^n = 0$$

and so the polynomial that fulfills  $f(a) = 0$  is  $f(x) = \sum_{i=0}^n f_i x^i$ . □

## 8.3 Polynomials

**Definition 8.7.** Suppose that  $K$  is a field extension of  $F$ . Let  $a \in K$  be algebraic over  $F$ , and let  $n \geq 0$  be the least integer such that there exists  $f(x) \in F[x]$  of degree  $n$  s.t.  $f(a)=0$ . Then  $n$  is called **the degree of  $a$** , and we mark

$$\mathit{deg}_F(a) = n$$

It's important to note that  $f$  from the above definition isn't unique (It's actually unique if you include multiplication with  $f \in F$ , and also if the definition states that  $f$  is monic, meaning  $f_0 = 1$ ).

The degree of  $a$ ,  $n = \mathit{deg}_F(a)$  is, however, unique.

**Definition 8.8.** The monic polynomial that applies to the prior definition is called **the minimal polynomial**

**Claim 9.** *The minimal polynomial is irreducible over  $F[x]$ .*

*Proof.* Assume by way of contradiction that  $p(x)$  is the minimal polynomial and that it is reducible, meaning it can be written as  $p(x) = f(x)g(x)$ , where  $f(x), g(x)$  are irreducible with  $\mathit{deg} g \geq 1, \mathit{deg} f \geq 1$ .  $a$  is algebraic over  $F$  for that  $f(x)$ , and so  $p(a) = f(a) \cdot g(a) = 0$ .  $F$  is an integral domain and so the multiplication equals zero implies that either  $f(a) = 0$  or  $g(a) = 0$ . This means that one of those polynomials is the minimal polynomial for  $a$ , a contradiction.  $\square$

**Definition 8.10.** Suppose that  $K$  is a field extension of  $F$ . Let  $a \in K$  be algebraic over  $F$ . Let

$$F(a) = \{f(a) \mid f(x) \in F[x]\} \subset K$$

(This actually means to assign  $a$  (algebraic over  $F$ ) in all polynomials of  $F$ , and group them all together).

Our goal now will be to prove that  $F(a)$  is also a field.

**Theorem 8.11.**  *$F(a)$  is a field.*

Before we prove the above Theorem, let's present a helpful lemma.

**Lemma 8.12.** Assume that  $D, F$  are fields,  $D$  is a finite field extension of  $F$ . In addition,  $D$  is an integral domain. So  $D$  is a field.

*Proof.* Take  $0 \neq d \in D$ . To show that  $D$  is a field, since it's an i.d. all we need to show is that  $d$  has an inverse,  $d' \in D$ .

We know that  $D$  is a finite extension, namely:  $[D:F] = n < \infty$ .

That means, as we've seen, that  $D$  is algebraic over  $F$ , hence  $d$  is algebraic over  $F$ .

Based on that, we know that there exists  $f(x) \in F(x)$  such that  $f(d) = 0$ .

In other words, there is a polynomial  $f(x)$  s.t.

$$f_0 + f_1d + f_2d^2 + \dots + f_nd^n = 0$$

There are now two options: In case that  $f_0 \neq 0$  :

$$f(d) = f_1d + f_2d^2 + \dots + f_nd^n = -f_0 \implies \frac{-f_1}{f_0}d + \frac{-f_2}{f_0}d^2 + \dots + \frac{-f_n}{f_0}d^n = 1$$

Equivalently:

$$\underbrace{d\left[\frac{-f_1}{f_0}d^0 + \frac{-f_2}{f_0}d^1 + \dots + \frac{-f_n}{f_0}d^{n-1}\right]}_{\in D} = 1$$

Hence - We've found  $d$ 's inverse (in brackets)!

In case that  $f_0 = 0$ :

$$d[f_1 + f_2d + \dots + f_nd^{n-1}] = 0$$

Since  $D$  is an integral domain and  $d \neq 0$ , we can conclude that the the element inside the brackets equals 0.

By induction, it is easy to see that  $d$  has an inverse.

All in all, we see that every  $d \in D$  has an inverse in  $D$ , hence  $D$  is a field. □

Now, let's prove the above theorem:

*Proof.* In order to use the lemma and deduce that  $F(a)$  is a field, all we need is to prove that  $F(a)$  is a finite field extension of  $F$ .

$a$  is algebraic  $\implies$  there exists  $p(x) \in F[x]$  such that  $p(a) = 0$ .

mark  $\deg(p) = \deg_F(a) = n$ .

Let  $f(x)$  some polynomial  $F[x]$ . Looking at  $f(x) - r(x)p(x)$  - the polynomial that creates  $r(x)$ . We can write:

$$f(x) = g(x)p(x) + r(x)$$

when either  $r(x) = 0$  or  $r(x) \neq 0$  and  $\deg(r) < \deg(p)$ .

$f(x) = g(x)p(x) + r(x)$  and  $p(x) = 0$ , therefore  $f(x) = r(x)$ , and we derive that for any  $x$  we can find a polynomial of a lesser degree than  $f(x)$ .

We also know that for every  $f$ , the maximal degree is  $n$ . Thus:

$$F(a) = \{f(x) \mid f(x) \in F[x]\} = \{f(x) \mid f(x) \in F[x] \wedge \deg(f) \leq n\} = \text{span}_F(1, a, a^2, \dots, a^{n-1})$$

$$\Rightarrow [F(a) : F] \leq n$$

We have shown that  $F(a)$  is a finite field extension of  $F$ . Now, using the above lemma, it is immediate that  $F(a)$  is a field.  $\square$

Recap: Let  $K$  be a field extension of  $F$  and  $a \in K$  s.t.  $\deg(a) = n$ .

Then, as we've seen:  $F(a)$  is a field extension of  $F$  and  $[F(a) : F] \leq n$ . We now want to show that it is exactly  $n$ .

**Claim 13.**  $[F(a) : F] = n$ .

*Proof.* Showing a set of  $n$  elements in  $F(a)$  that are linearly independent will prove the claim. Let's take the following set:

$$\{1, a, a^2, \dots, a^{n-1}\}$$

Assume that this set's elements are linearly dependent over  $F$ : there exist  $f_0, f_1, f_2, \dots, f_{n-1} \in F$ , not all zeros, such that

$$f_0 + f_1 a + f_2 a^2 + \dots + f_{n-1} a^{n-1} = 0$$

This is a polynomial of  $\deg n-1$  s.t.  $f(a) = 0$ . This is a contradiction to the fact that  $a$  is algebraic over  $F$  with  $\deg a = n$ . So our assumption must be false, and these  $n$  elements are linearly independent, as required.  $\square$



### 8.3.1 Roots of polynomials

**Definition 8.14.** (Root of polynomial) Let  $F \subseteq K$  and  $K$  is a field extension of  $F$ .  $a \in K$  is a **root** of  $f(x) \in F[x]$  if  $f(a) = 0$ .

**Claim 15.**  $a$  is a **root** of  $f(x) \iff x - a \mid f(x)$  in  $K[x]$ .

*Proof.*  $\leftarrow$ : Let  $x - a \mid f(x)$  then exists  $g(x) \in K[x]$  s.t.  $f(x) = (x - a) \cdot g(x)$ , it easy to see that  $f(a) = (a - a) \cdot g(a) = 0$ .

$\rightarrow$ : Let  $f(a) = 0$  and let write  $f(x)$  as  $f(x) = (x - a) \cdot g(x) + r(x)$  for some polynomials  $g(x)$  and  $r(x)$  in  $K[x]$ . we know from the division of  $f(x)$  in  $(x - a)$  that  $\deg(r(x)) < 1$  (since  $\deg(x - a) < 1$ )  $\Rightarrow r(x)$  is a constant. After putting all together we get  $0 = f(a) = (a - a) \cdot g(a) + r(a) = r(a) \Rightarrow r(a) = 0 \Rightarrow r(x) = 0 \Rightarrow f(x) = (x - a) \cdot g(x)$ .  $\square$

**Corollary 8.16.** Let  $f(x) \in F[x]$  with  $\deg(f(x)) = n$ ,  $f(x)$  can have at most  $n$  roots in any field extension of  $F$ .

*Proof.* By induction on  $n$  above any field.

•  $n=1$ , trivial. • Let assume that every  $g(x) \in F[x]$  with  $\deg(g(x)) = n - 1$  has at most  $n-1$  roots in any field extension of  $F$ . • Step: Let  $f(x) \in F[x]$ , with  $\deg(f(x)) = n$  and  $a$  is a root of  $f(x)$ , then  $f(x)$  can be written as:

$f(x) = (x - a) \cdot g(x)$  for  $g(x) \in K[x]$  field extension of  $F$ . Since  $\deg(x - a) = 1$ , and the isomorphism between product of polynomials to sum of degrees,  $\deg(g(x)) = n - 1$  for any  $g(x) \in K[x]$  that satisfies the equation above. From the assumption,  $g(x)$  has at most  $n-1$  roots, and including  $a$ ,  $f(x)$  has at most  $n$  roots.  $\square$

**Claim 17.** Let  $f(x) \in F[x]$  Then there exist  $K[x]$  field extension of  $F$  s.t.

$f(x) = \prod_{i=1}^n (x - a_i)$ . in other words, for any polynomial  $f(x)$  with degree  $n$ , exist a field extension that we can find all  $n$  roots of  $f(x)$ . For example, the roots of  $x^2 + 1 \in \mathbf{R}[x]$ , can be found in  $\mathbf{C}[x]$  that is a field extension of  $\mathbf{R}[x]$ .

*Proof.* Let  $f(x) = p(x) \cdot q(x)$ ,  $p(x) \in F[x]$  and  $p(x)$  is irreducible w.r.t.  $F[x]$ . Then,  $p(t) \in (F[x] / \langle p(x) \rangle)[t] \Rightarrow p(x) = 0 \Rightarrow t - x \mid p(t)$  (from the claim above)  $\Rightarrow p(t) = (t - x) \cdot g(t)$ .  $\deg(g(t)) = \deg(p(t)) - 1$ . We continue this process for all the factors of  $f(x)$ .  $\square$

### 8.3.2 Derivatives: formal derivative

**Definition 8.18.**  $f(x) \in F[x]$  :

- $(x^n)' = n \cdot x^{n-1}$
- $(f + g)' = f' + g'$
- $(c \cdot f(x))' = c \cdot f'(x)$

*Property 8.18.1.*  $(f \cdot g)' = f' \cdot g + f \cdot g'$

**Claim 19.** Let  $f(x) \in F[x]$ , without know any of its roots, then  $f(x)$  has a repeted root (i.e,  $\exists a \in K$  s.t  $(x - a)^2 | f(x)$ )  $\iff \langle f \rangle + \langle f' \rangle \neq \langle 1 \rangle$ . In less formal words:  $\gcd(f(x), f'(x)) \neq 1$ .

*Proof.*  $\rightarrow f(x) = (x - a)^2 \cdot g(x) \Rightarrow f'(x) = 2 \cdot (x - a) \cdot g(x) + (x - a)^2 \cdot g'(x)$   
with  $(x - a)^2 \in K, g(x) \in K[x]$ .  $\Rightarrow (x - a) | f'(x), (x - a) | f(x) \Rightarrow \langle f(x) \rangle + \langle f'(x) \rangle \subseteq \langle x - a \rangle \neq \langle 1 \rangle$

Second direction is left as homework. □

## LECTURE 9

# FINITE FIELDS CONT.; SMALL-BIASED SETS

---

In this lecture we will finish discussing finite fields, present a construction for  $\mathbb{F}_{p^n}$  and present the problem of small-biased sets and some of the proposed solutions to it.

## 9.1 Extension Fields

### 9.1.1 Recap

In the last chapter, we have seen that by taking some  $f(x) \in F[x]$  for some field  $F$ , we may not have any or all roots of  $f(x)$  in  $F$  itself. However, we can extend  $F$  to a field  $Q$  in which  $f$  will have all its roots.

We have already shown that if  $f(x)$  is irreducible, then  $K := F[x]/\langle f(x) \rangle$  is a field, in which we view  $f(x) \in F[x]$  as  $f(y) \in K[y]$  where  $f(y) = (y-a)g(a)$ ,  $a \in K$ ,  $g(y) \in K[y]$ . We can continue this process until we find all roots of  $f$ , as every time we perform this procedure, the degree decreases ( $\deg(g) < \deg(f)$ ).

### 9.1.2 Splitting Fields

**Theorem 9.1.**  $\forall f(x) \in F[x]$  where  $f(x)$  is not constant, there exists an extension field  $K$  of  $F$  (which depends both on  $F$  and  $f$ ) so that  $f(x)$  has all of its roots in  $K$  and in particular  $f(x) = \prod_{i=1}^{\deg(f)} (x - a_i)$ ,  $a_i \in K$  (where there may be some  $i \neq j$  with  $a_i = a_j$ )

*In addition,  $\forall F \subseteq K$ ,  $a_1, \dots, a_n \in K$ , The smallest field which contains  $F$  and  $a_1, \dots, a_n$  is unique up to isomorphism. It is called the **Splitting Field of  $f$  over  $F$***

### 9.1.3 $F(a_1, \dots, a_n)$

**Definition 9.2.**  $F(a)$  all the rational functions of  $a$  with coefficients from  $F$

**Corollary 9.3.**  $F(a_1, \dots, a_n)$  is all the rational functions of  $a_1, \dots, a_n$  with coefficients from  $F$  which can be seen recursively as all the rational functions over  $a_2, \dots, a_n$  with coefficients from  $F(a_1)$ .  $F(a_1, \dots, a_n) = (F(a_1))(a_2, \dots, a_n)$

*Example 9.4.*  $\frac{ab^2}{a+1} + \frac{b}{a^3} \in F(a, b)$  where the variables are  $a, b$  and the coefficients are all  $1 \in F$ . But this can also be viewed as  $(\frac{a}{a+1})b^2 + (\frac{1}{a^3})b \in (F(a))(b)$  where the variable is  $b$  and the coefficients are  $(\frac{a}{a+1}), (\frac{1}{a^3}) \in F(a)$

## 9.2 Finite Fields

Previously, we have shown that for any primary number  $p \in \mathbb{N}$  there exists a field  $\mathbb{F}_p$  where addition and multiplication are done as in  $\mathbb{Z} \pmod p$

### 9.2.1 Existence of $\mathbb{F}_{p^n}$

*Reminder 1* (Fermat's little theorem). If  $p$  is primary, then for any  $x \in \mathbb{Z}$  it holds that  $x^p \equiv x \pmod p$

**Claim 5** (Existence of  $\mathbb{F}_{p^n}$ ). Let there be some  $p \in \mathbb{N}$  primary, and  $f(x) \in \mathbb{F}_p$  irreducible with  $\deg(f) = n$  (assuming there is such  $f(x)$ ), then  $\mathbb{F}_p[x]/\langle f(x) \rangle$  is a field of size  $p^n$ . I.e.  $[K : \mathbb{F}_p] = n \Rightarrow |K| = p^n$ .

**Corollary 9.6** ( $K$  as a vector space over  $\mathbb{F}_p$ ). There are  $v_1, \dots, v_n$  a base for  $K$  over  $\mathbb{F}_p$

*Reminder 2.*  $K$  is a vector space over  $\mathbb{F}_p$  with base  $v_1, \dots, v_n \in K$  if  $\forall k \in K, k$  can be written uniquely as  $\sum_{i=1}^n a_i v_i$  where  $a_1, \dots, a_n \in \mathbb{F}_p$

*Remark.* An extension field is always a vector space over the base field, thus its order is a power of the order of the base field.

## 9.2.2 Construction of $\mathbb{F}_{p^n}$

Let  $K$  be a field s.t.  $|K| = p^n$ , then  $\forall x \in K \ x^{|K|} = x^{p^n} = x$ .

This is true because  $K \setminus \{0\}$  is a multiplicative group of size  $p^n - 1$ , and therefore  $x^{p^n} = x^{p^n-1}x = x$ . This means that in a sense, the polynomial  $x^{|K|} - x \in \mathbb{F}_p[x]$  holds all the information about  $K$  (its roots are exactly the elements of  $K$ ).

Consider the Polynomial  $x^{p^n} - x \in \mathbb{F}_p[x]$ , where  $p$  is prime and  $n \geq 1$ . We can use this polynomial to construct an extension of  $\mathbb{F}_p$  with size  $p^n$ . We know that there exists some extension of  $\mathbb{F}_p[x]$ ,  $K$ , s.t. all the roots of  $x^{p^n} - x$  are in  $K$ .

Denote the roots of this polynomial by  $R = \{a_1, a_2, \dots, a_{p^n}\}$ , note that there could be repetitions in  $R$ , but we will show later that there aren't.

**Claim 7.**  $R$  is a field

To show that  $R$  is a field we need to prove that  $\forall a, b \in R$ :

1.  $ab \in R$
2.  $a + b \in R$
3.  $a^{-1} \in R, a \neq 0$

The other required properties are trivial.

proof:

$$1. \ a^{p^n} = a, \ b^{p^n} = b \implies ab = a^{p^n}b^{p^n} = (ab)^{p^n} \implies ab \in R$$

$$2. \ (a + b)^{p^n} = a^{p^n} + \binom{p^n}{1}a^{p^n-1}b + \binom{p^n}{2}a^{p^n-2}b^2 + \dots + b^{p^n}$$

Notice that for all  $1 \leq i \leq p^n - 1$ ,  $\binom{p^n}{i} = \frac{p^n!}{i!(p^n-i)!} = p^n \cdot (\text{something}) = 0 \pmod{p^n}$   
 (since the characteristic of this field is  $p^n$ , therefore we get  $(a+b)^{p^n} = a^{p^n} + b^{p^n} = a + b$ , i.e.  $a + b \in R$ )

$$3. \ a \neq 0, \ (a^{-1})^{p^n} = (a^{p^n})^{-1} = a^{-1} \implies a^{-1} \in R$$

□

It is left to show that there are no repetitions in  $R$ . We know from the previous lecture, that a polynomial  $f$  has repetitions in its roots iff  $\langle f \rangle + \langle f' \rangle \neq \langle 1 \rangle$ . In

our case,  $(x^{p^n} - x)' = p^n x^{p^n-1} - 1 = -1 \pmod{p}$ , hence  $\langle x^{p^n} - x \rangle + \langle (x^{p^n} - x)' \rangle = \langle x^{p^n} - x \rangle + \langle -1 \rangle = \langle x^{p^n} - x \rangle + \langle 1 \rangle = \langle 1 \rangle$ , therefore there are no repetitions in  $R$ . It is easy to confirm that  $\mathbb{F}_p \subseteq R$  by assigning any element from  $\mathbb{F}_p$  to  $x^{p^n} - x$ . Finally we get that  $R$  is an extension of  $\mathbb{F}_p$  with size  $p^n$ .

### 9.3 Small Biased Sets

**Definition 9.8** ( $\epsilon$ -biased set).  $s \in \{0, 1\}^n$  is called  $\epsilon$ -biased set if

$$\forall \tau \in \{0, 1\}^n \setminus \{0\} : \left| \mathbb{E}_{s \sim S} [(-1)^{\langle s, \tau \rangle}] \right| \leq \epsilon$$

Trying to phrase the definition in a more intuitive manner: Let's randomly select an element from  $S$  and look at its inner product with some vector  $\tau: \sum_{i=1}^n s_i \tau_i$ . We want the probability to get an even result to be approximately equal (up to  $\epsilon$ ) to the probability to get an odd result.

**Definition 9.9** (Pseudorandom Distribution). Let  $C$  be a set of functions:  $C = \{f : \{0, 1\}^n \rightarrow \{0, 1\}\}$ . We want to create a distribution  $D$  over  $\{0, 1\}^n$ , such that no member of  $C$  will be able to separate between a sample from  $D$  and a sample from the uniform distribution. Formally:  $D$  will be called  $\epsilon$ -Pseudorandom w.r.t  $C$ , if:

$$\forall f \in C : |Pr[f(u_n) = 1] - Pr[f(d_n) = 1]| \leq \epsilon$$

( $u_n$  is a sample taken from the uniform dist.,  $d_n$  is a sample taken from  $D$ ).

We shall notice that switching the 1 in the definition above into a 0 will create a logically equal definition. With this insight, we can say that a distribution  $D$  is pseudorandom iff:

$$\forall f \in C : |\mathbb{E}[f(u_n)] - \mathbb{E}[f(d_n)]| \leq \epsilon$$

Goal: given  $n, \epsilon$  we want to build an  $n, \epsilon$  biased group  $S \subseteq \{0, 1\}^n$  which is smallest possible. Our set  $S$  induces distribution  $D$  because we sample uniformly at random from the set. Namely, no linear test could distinguish between sample from  $U$  and sample from  $D$ , where  $U$  is the distribution which sample uniformly from  $\{0, 1\}^n$ .

### 9.3.1 Remarkable results

It was proven that a tight bound on size of the set exists, but is computationally expensive. That bound is:  $\frac{n}{\epsilon^2}$ .

We shall now mention other remarkable results in the field of small-biased sets, in which researchers were able to find sets of small size:

Publisher	Set Size
Naor and Naor 1980	$\frac{n}{\epsilon^{\Theta(1)}}$
AGHP 1992	$\frac{n^2}{\epsilon^2}$
Amnon Ta Shma and Avi Ben Aroya	$(\frac{n}{\epsilon^2})^{\frac{5}{4}}$
Amnon Ta Shma	$\frac{n^2}{\epsilon^{(2+\Theta(1))}}$

### 9.3.2 The powery construction (AGHP)

In this section we will create an  $\epsilon$  small-biased set of size  $\Theta((\frac{\epsilon}{n})^2)$ . To do so, will use the notation  $l \in \mathbb{N}$ , and will only commit to a value for said notation near the end of the section.

Now - let's build our set. Let  $\mathbb{F}_{2^l}$  be a finite field (we've seen its existence in the previous section).

$\forall x, y \in \mathbb{F}_{2^l}$  define  $S_{xy} \in \{0, 1\}^n : (S_{xy})_i = \langle x^i, y \rangle_{i \in \{1 \dots n\}}$

Therefore:

$$(S_{xy}) = (\langle x^1, y \rangle, \langle x^2, y \rangle, \dots, \langle x^n, y \rangle)$$

Our set shall be:  $S = \{S_{xy} | x, y \in \mathbb{F}_{2^l}\}$ , and its size:  $|S| = m = 2^{2l} = 2^{2l}$

In the previous definition, we have used the  $\langle, \rangle$  notation - the inner product of a vector space. Since  $\mathbb{F}_{2^l}$  is a field - how is this legal? Lucky for us  $\mathbb{F}_{2^l} \simeq \mathbb{F}_2^l$  where  $\mathbb{F}_2^l$  is a vector space over  $\mathbb{F}_2$ . This allow us to use the inner product of  $\mathbb{F}_2^l$ .

As mentioned before - we would like to show that S is a small-biased set. To do so, let's look at the inner product of a vector  $S_{xy} \in S$  and  $0 \neq \tau \in \{0, 1\}^n$ :

$$\langle S_{xy}, \tau \rangle = \sum_{i=1}^n ((S_{xy})_i \cdot \tau_i) = \sum_{i=1}^n (\langle x^i, y \rangle \cdot \tau_i) = \sum_{i=1}^n \langle \tau_i \cdot x^i, y \rangle = \langle \sum_{i=1}^n \tau_i \cdot x^i, y \rangle.$$

Let us notice that  $\forall \tau \in \{0, 1\}^n$ , we can define a polynomial  $P_\tau(x) = \sum_{i=1}^n \tau_i \cdot x^i$ . Using this notation:  $\langle S_{xy}, \tau \rangle = \langle P_\tau(x), y \rangle$ .

If  $x$  is a root of  $P_\tau(x)$  then we know that  $\langle P_\tau(x), y \rangle = \langle \bar{0}, y \rangle = 0$ . In any other case (since  $0 \neq \tau$ ), we have that  $P_\tau(x) \neq 0$  which means:

$$\mathbb{E}_{x, y \sim \mathbb{F}_{2^l}} [(-1)^{\langle S_{xy}, \tau \rangle}] = 0$$

From all of the above, we can understand that:

$$\begin{aligned} \mathbb{E}_{x, y \sim \mathbb{F}_{2^l}} [(-1)^{\langle S_{xy}, \tau \rangle}] &= \mathbb{E}_{y \sim \mathbb{F}_{2^l}} [ \mathbb{E}_{x \sim \mathbb{F}_{2^l}} [(-1)^{\langle S_{xy}, \tau \rangle}] | P_\tau(x) = 0 ] \cdot Pr(P_\tau(x) = 0) + \\ &\quad \mathbb{E}_{y \sim \mathbb{F}_{2^l}} [ \mathbb{E}_{x \sim \mathbb{F}_{2^l}} [(-1)^{\langle S_{xy}, \tau \rangle}] | P_\tau(x) \neq 0 ] \cdot Pr(P_\tau(x) \neq 0) = \\ &\quad Pr(P_\tau(x) = 0) + 0 \end{aligned}$$

In conclusion, we get that:

$$\mathbb{E}_{x, y \sim \mathbb{F}_{2^l}} [(-1)^{\langle S_{xy}, \tau \rangle}] = Pr(P_\tau(x) = 0) \leq \frac{n}{2^l}$$

It is finally time to choose  $l$ . We shall the smallest  $l$  possible such that:  $\frac{n}{2^l} \leq \epsilon$ .

This means that we have found  $S$  - an  $\epsilon$  small-biased set of size  $|S| = m = 2^{2l} = \Theta\left(\left(\frac{n}{\epsilon}\right)^2\right)$ , just like we wanted.



# LECTURE 10

## SMALL BIAS SETS

---

We'll see an explicit construction of an  $\epsilon$ -biased set over  $k$  bits of size  $O\left(\frac{n}{\epsilon^2}\right)^{\frac{5}{4}}$ .

### 10.1 Bezout Theorem

**Definition 10.1.** A multivariate polynomial  $f \in \mathbb{F}[x_1, \dots, x_n]$  has a degree,  $d$ , defined by:

$$d := \text{TotalDeg}(f) = \deg(f(x, x, \dots, x))$$

**Theorem 10.2.** (*Bezout Theorem*) For any  $f, g \in \mathbb{F}[x, y]$  and for:

$$I = \{(x, y) \in \mathbb{F}^2 : f(x, y) = g(x, y) = 0\}$$

assuming  $\langle f \rangle + \langle g \rangle = \langle 1 \rangle$  (as ideals in  $\mathbb{F}[x, y]$ ):

$$|I| \leq \text{deg}(f) \cdot \text{deg}(g)$$

*Remark.* The assumption  $\langle f \rangle + \langle g \rangle = \langle 1 \rangle$  is necessary: if we assume  $\langle f \rangle + \langle g \rangle \neq \langle 1 \rangle$  (generalization of GCD with ideals, meaning there's a non-trivial polynomial that divides both), if there's an  $h$  s.t  $f = h \cdot a$  and  $g = h \cdot b$ , then  $f$  and  $g$  share all the zeros of  $h$ . For example in the case  $h(x, y) = x^2 + y^2 - 1$ , it implies  $|I| = \infty$  because  $h$  has infinite zeros.

*Example 10.3.* Given a parabola  $g = y - x^2$  and a circle  $f = x^2 + y^2 - 1$ , by Bezout theorem  $|I| \leq 4$ . There might be 4 or 2 intersections.

### 10.2 Ben-Aroya, Ta-Shma Construction

#### 10.2.1 The construction

**Definition 10.4.** (Hermitian Curve) Let  $p = 2^l$  and  $q := p^2$  for some  $l \in \mathbb{N}$ . Consider the equality  $y^p + y = x^{p+1}$  over  $\mathbb{F}_q = \mathbb{F}_{p^2}$ . The Hermitian Curve is defined

as:

$$H = \{(x, y) \in \mathbb{F}_q^2 : y^p + y - x^{p+1} = 0\}$$

Given  $p, q, l, H$  as defined above, construct  $S$  as follows

$$S = \{s(x, y, z) : (x, y) \in H, z \in \mathbb{F}_q\}$$

where  $s(x, y, z)_{ij} = \langle x^i, y^j, z \rangle$ .

We choose indices  $i, j$  so there are  $n$  entries in the vector  $s(x, y, z)$ :  $i = 1, \dots, r$  and  $j = 1, \dots, \frac{n}{r}$ . So  $i, j \in \{(i, j) : i + j \leq \sqrt{2n}\}$ .

## 10.2.2 Analyzing the Construction

*Remark.* We'll show that  $\forall \tau \in \{0, 1\}^n \setminus \{0^n\}$ ,  $|\mathbb{E}[(-1)^{\langle s, \tau \rangle}]| \leq \epsilon$ . Last lecture we've defined  $\langle s, \tau \rangle = \sum_{i=1}^n \langle x^i, y \rangle \tau_i$  and showed that for  $p_\tau(x) = \sum_{i=1}^n \tau_i x^i$  the term equals to  $\langle p_\tau(x), y \rangle$ .

We'd like to show that:

$$\left| \mathbb{E}_{\substack{(x,y) \sim H \\ z \sim \mathbb{F}_q}} [(-1)^{\langle s(x,y,z), \tau \rangle}] \right| = \left| \mathbb{E}_{\substack{(x,y) \sim H \\ z \sim \mathbb{F}_q}} [(-1)^{\sum_{i+j \leq \sqrt{2n}} s(x,y,z)_{ij} \tau_{ij}}] \right| \leq \epsilon$$

As in the last construction, we'll notice we can include  $\tau_{ij}$  in the product:

$$\sum_{i+j \leq \sqrt{2n}} \langle x^i y^j, z \rangle \tau_{ij} = \sum_{i+j \leq \sqrt{2n}} \langle \tau_{ij} \cdot x^i y^j, z \rangle = \left\langle \sum_{i+j \leq \sqrt{2n}} \tau_{ij} \cdot x^i y^j, z \right\rangle = \langle f_\tau(x, y), z \rangle$$

Where  $f_\tau(x, y) = \sum_{i+j \leq \sqrt{2n}} \tau_{ij} \cdot x^i y^j$ . Hence we have:

$$\mathbb{E}_{\substack{(x,y) \sim H \\ z \sim \mathbb{F}_q}} [(-1)^{\sum_{i+j \leq \sqrt{2n}} s(x,y,z)_{ij} \tau_{ij}}] = \mathbb{E}_{\substack{(x,y) \sim H \\ z \sim \mathbb{F}_q}} [(-1)^{\langle f_\tau(x,y), z \rangle}]$$

Notice that we transferred a "combinatorial game" to a "polynomials game". The question we're asking is, how many zeros  $f_\tau(x, y)$  has in  $H$ ? Specifically, we're interested in the ratio between this number and the size of  $H$ , meaning  $\frac{\#\text{ roots of } f_\tau \text{ in } H}{|H|}$  (recall we take the expectation over samples  $(x, y) \in H$ ).

*Remark.* If  $f_\tau(x, y) \neq 0$ , by choosing  $z \in \mathbb{F}_q$  randomly the expectation is 0:

$$\mathbb{E}_{z \sim \mathbb{F}_q} [(-1)^{\langle f_\tau(x, y), z \rangle}] = 0$$

Hence, we only care about roots of  $f_\tau$  in  $H$ .

Notice that  $\deg(f) \cdot \deg(H) \leq \sqrt{2n} \cdot (p+1) = O(p\sqrt{n})$ . By Bezout Theorem, this implies the number of  $f_\tau$  roots in  $H$  is bounded by  $O(p\sqrt{n})$  (we'll prove the condition  $\langle f \rangle + \langle H \rangle = \langle 1 \rangle$  later). It is left to show that  $|H|$  is big enough.

**Claim 5.**  $|H| = p^3 = pq$

**Claim 6.** *The polynomial  $H(x, y) = y^p + y - x^{p+1}$  is irreducible as an element of  $\mathbb{F}_q[x, y]$ .*

**Corollary 10.7.** *Given the claims, it follows that  $\langle f_\tau \rangle + \langle H \rangle = \langle 1 \rangle$ .*

*Proof.* Choose  $l = \log n$ , then  $\deg(H) = n+1$  (recall that  $p = 2^l$ ), and  $H$  is irreducible. Assuming towards contradiction that  $\langle f_\tau \rangle + \langle H \rangle \neq \langle 1 \rangle$ , there is some polynomial  $h \in \mathbb{F}[x, y]$  s.t.  $h|f_\tau$  and  $h|H$ ; as  $H$  is irreducible and  $h$  is non-trivial we must have  $h = H$ , which implies that  $H|f_\tau$ . Then  $\deg(H) = n+1 > \sqrt{2n} \geq \deg(\mathbb{F}_\tau)$  in contradiction.  $\square$

Therefore, by Bezout theorem we get the bound:

$$|\mathbb{E}[(-1)^{\langle s(x, y, z), \tau \rangle}]| = \frac{\# \text{ roots of } f_\tau \text{ in } H}{|H|} \leq \frac{p\sqrt{n}}{p^3} = \frac{\sqrt{n}}{p}$$

*Remark.* It is enough to choose  $l$  s.t  $p = 2^l \gg \sqrt{n}$ .

We get  $\epsilon = \frac{\sqrt{n}}{p^2}$  and so  $p^4 = \frac{n}{\epsilon^2}$ . Notice that

$$|S| = |H| \cdot |\mathbb{F}_q| = p^3 \cdot p^2 = p^5$$

With  $\epsilon$  above we get

$$|S| = (n/\epsilon^2)^{\frac{5}{4}}$$

### 10.2.3 Missing Proofs

*Proof.* (Claim 1) It is enough to show that  $y^p + y - x^{p+1} = 0$  has exactly  $pq$  solutions. Denote  $f(y) = y^p + y$  and  $H$  the hermetian curve as before.

Consider  $f$  as a transformation  $f : \mathbb{F}_q \rightarrow \mathbb{F}_q$ , then  $Im(f) \subseteq \mathbb{F}_p \subset \mathbb{F}_q$ .

Explanation: plug-in  $y^p + y$  in  $z^p - z$ . If the result is 0, it implies that  $y^p + y \in \mathbb{F}_p$  (the roots of  $z^p - z$  are in  $\mathbb{F}_p$ ). We have that

$$(y^p + y)^p - (y^p + y) = y^{p^2} + y^p - y^p - y = y^q - y = 0$$

Where the first equality follows from the fact that  $p = 2^l$  over a field with  $Char\mathbb{F} = 2$ . The second follows from the fact that  $p^2 = q$ , and the third follows from  $y \in \mathbb{F}_q$ . Meaning,  $y^p + y \in \mathbb{F}_p$  for any  $y \in \mathbb{F}_q$ , hence  $Im(f) \subseteq \mathbb{F}_p$ . Similarly, if  $x \in \mathbb{F}_q$  then  $x^{p+1} \in \mathbb{F}_p$ :

$$(x^{p+1})^p - x^{p+1} = x^{p^2} \cdot x^p - x^{p+1} = x^q \cdot x^p - x^{p+1} = x \cdot x^p - x^{p+1} = 0$$

For  $x = 0$  and  $y \in \mathbb{F}_p$ :

$$x^{p+1} = 0 = 2y = y + y = y^p + y$$

Meaning,

$$\{0\} \times \mathbb{F}_p = \{(0, y) : y \in \mathbb{F}_p\} \subseteq H$$

We'll notice that  $y^p + y = 0$  is a  $\mathbb{F}_p$ -linear function.

**Definition 10.8.** We'll say that a function  $f : \mathbb{F}_p^2 \rightarrow \mathbb{F}_p^2$  is a  $\mathbb{F}_p$ -linear function if  $\forall x, y \in \mathbb{F}_p^2$  and  $\forall a, b \in \mathbb{F}_p$ :

$$f(ax + by) = af(x) + bf(y)$$

It's not hard to show that  $f(t) = t^p + t$  is indeed  $\mathbb{F}_p$ -linear function. That is because  $y \rightarrow y$  and  $y \rightarrow y^p$  are a  $\mathbb{F}_p$ -linear functions and their sum is  $\mathbb{F}_p$ -linear too. Meaning  $y \rightarrow y^p + y$  is a linear function above  $\mathbb{F}_p$  so that the solutions space dimension is one less than the vector space  $\mathbb{F}_p$ ; meaning  $p^{2-1} = p$ .

We must show that  $\forall x \in \mathbb{F}_q$  the equation  $y + y^p = x^{p+1}$  has  $p$  solutions and for that we'll use a 'trick'.

$y + y^p \in \mathbb{F}_p$  for all  $y \in \mathbb{F}_q$ . If there's a solution  $y^p + y = x^{p+1}$  for a constant  $x \in \mathbb{F}_q$  there are exactly  $p$  solutions since  $y^p + y$  is linear.

There are  $p^2$  possible solutions for  $y^p + y$  (since there are  $p^2$  solutions for choosing  $y$ ). We get every result exactly  $p$  times. Therefore, for every  $a \in \mathbb{F}_p$  there are exactly  $p$  points  $y \in \mathbb{F}_q$  such that  $y^p + y = a$ . For every  $x \in \mathbb{F}_q$  it holds that  $x^{p+1} \in \mathbb{F}_p$  and therefore every  $x \in \mathbb{F}_q$  there are exactly  $p$  elements  $y \in \mathbb{F}_q$  such that  $H(x, y) = 0$ , meaning:

$$|H| = pq = p^3$$

□

For the proof of claim 2 we'll need to use Eisenstein criterion.

**Theorem 10.9.** (*Eisenstein criterion*)

Let there be  $R$  an integral domain and  $f(x) \in R[x]$ . Denote  $f(x) = \sum_{i=0}^n a_i x^i$  and assume that there exists a prime ideal  $P \subset R$  such that:

1.  $a_0, \dots, a_{n-1} \in P$
2.  $a_n \notin P$
3.  $a_0 \notin P^2$  (where  $P^2 = \{\sum_{j=0}^n p_{0,j} p_{1,j} : p_{0,j}, p_{1,j} \in P, n \in \mathbb{N}\}$ )

Under these assumptions,  $f$  is irreducible in  $R[x]$ .

*Proof.* (Eisenstein criterion's, from Wikipedia)

We'll assume  $f$  is reducible, i.e.  $f(x) = b(x)c(x)$  for some non-trivial  $b(x) = \sum_{i=0}^r b_i x^i$ ,  $c(x) = \sum_{i=0}^s c_i x^i$ .

It's clear that  $a_0 = b_0 c_0 \in P \setminus P^2$  and therefore it cannot be that both of  $b_0, c_0$  are in  $P$ .

If  $b_0, c_0 \notin P$  since  $P$  is prime  $a_0 = b_0 c_0 \notin P$  in contradiction that multiplication of elements in  $P$  gives an element in  $P$ .

So only one of  $b_0, c_0$  is in  $P$ . WLOG assume that  $b_0 \in P$  and  $c_0 \notin P$ . We'll notice that:

$$a_1 = b_0 c_1 + b_1 c_0 \Rightarrow b_1 c_0 = a_1 - b_0 c_1 \in P$$

$(b_0, a_1 \in P)$  If  $b_1 \notin P$  then  $b_1 c_0 \in P$  for  $c_0, b_1 \notin P$  - in contradiction to the primness of  $P$ , so  $b_1 \in P$ .

Inductively, assuming  $\deg(c) \geq 1$ , assume  $b_0, b_1, \dots, b_{k-1} \in P$  for  $k \leq r < n$  ( $r < n$  because  $\deg(c) \geq 1$ ). As  $k < n$ ,  $a_k \in P$  and:

$$a_k = b_0c_k + b_1c_{k-1} + \dots + b_r c_0 = \sum_{i=0}^k b_i c_{k-i} \in P$$

where  $c_i = 0$  if  $i > s$ . Particularly:

$$b_r c_0 = a_k - \sum_{i=0}^{k-1} b_i c_{k-i} \in P$$

because  $a_k, b_0, b_1, \dots, b_{k-1} \in P$ . As  $c_0 \notin P$ , if  $b_k \notin P$  we contradict  $P$ 's primality again. Thus,  $b_0, b_1, \dots, b_r \in P$  by induction.

But then we get a contradiction:  $b_r c_{n-r} = a_n$  but  $a_n \notin P$  and  $b_r c_{n-r} \in P$ !

Therefore  $f$  is irreducible. □

*Proof.* (Claim 2)

To prove  $y^p + y - x^{p-1}$  is irreducible over  $\mathbb{F}_q[x, y]$ , we'll show that it's irreducible over  $\mathbb{F}_q[y][x]$  (which is the same).

We'll set  $R = \mathbb{F}_q[y]$  and denote  $P = \langle y \rangle \triangleleft R$ .

*Remark.*  $P$  is prime. That is because  $\mathbb{F}_q[y]/\langle y \rangle \cong \mathbb{F}_q$  and is therefore a field, so  $P$  is a maximal ideal and therefore is also prime (because  $\mathbb{F}_q[y]$  is PID).

It holds that:

$$R[x] \ni f_y(x) = -x^{p+1} + y^p + y = a_1x + a_0$$

where  $a_1 = -1$  and  $a_0 = y + y^p$  (note that  $a_0$  is an element of the ring  $R = \mathbb{F}_q[y]$ ). It's clear that  $a_0 \in P$ . What is  $P^2$ ? If  $\alpha \in P^2 = \langle y \rangle^2$  then there exist  $g_j, h_j \in \langle y \rangle$  s.t.  $\alpha = \sum_j g_j h_j$ ; but  $y|g_j$  and  $y|h_j$  for every  $j$ , so  $y^2 | \sum_j g_j h_j = \alpha$ , or  $P^2 \subseteq \langle y^2 \rangle$ . Obviously  $y^2 \in P^2$ , so  $P^2 = \langle y^2 \rangle$ . Thus:

1.  $a_0 = y(y^{p-1} + 1) \in P$
2.  $a_1 = a_n = -1 \notin P$  (because if  $-1 \in P$  then  $P = R$ )
3.  $a_0 \notin P^2 = \langle y^2 \rangle$  (because  $y^2 \nmid y^p + y$ )

Therefore,  $f_y \in R[x]$  satisfies the Eisenstein criterion's conditions. As such,  $f_y$  is irreducible in  $R[x]$  which immediately implies that  $f(x, y)$  is irreducible in  $\mathbb{F}_q[x, y]$ .

□

# LECTURE 11

## RANDOMNESS MERGERS

---

### 11.1 Small-Biased Sets

#### 11.1.1 Introduction

In the previous lectures we saw two constructions of small-biased sets:

- **AGHP**, which gives a small-biased set of size  $O\left(\frac{\epsilon}{n}\right)^2$ .
- **Ben Aroya - Ta-Shma**, which gives a small-biased set of size  $O\left(\frac{n}{\epsilon^2}\right)^{5/4}$ .

In this part, we will show that there exists a small-biased set of size  $O\left(\frac{n}{\epsilon^2}\right)$  using the probabilistic method. Unfortunately, the proof won't be constructive - it will show only the existence of such sets, without showing an explicit construction. First, we start by explaining the general approach of the probabilistic method, and then showing the case for small-biased sets.

#### 11.1.2 The Probabilistic Method

If every object in a collection of objects fails to have a certain property, then the probability that a random object chosen from the collection has that property is zero.

Similarly, showing that the probability is strictly less than 1 can be used to prove the existence of an object that does not satisfy the prescribed properties.

For instance in our case, given a sample space of sets, the property will be "the set is not small-biased". We will prove that the probability of a randomly selected set to satisfy said property is strictly less than 1 - meaning there must exist a small-biased set.



### 11.1.3 Existence proof for a small-biased set

**Theorem 11.1** (Existence of small-biased set). *There exists an  $\epsilon$ -biased set of size  $O\left(\frac{n}{\epsilon^2}\right)$ .*

*Proof.* Let  $\Omega$  be the sample space consisting of all sets of size  $m$  that are subsets of  $\{0, 1\}^n$ . We will show that there exists a set  $S \in \Omega$  that is  $\epsilon$ -biased.

Let  $V_1, \dots, V_m \in \{0, 1\}^n$  be uniformly randomly chosen and independent. Let

$$S = \{V_1, \dots, V_m\}$$

We allow repetitions. Fix some linear test  $\tau \in \{0, 1\}^n \setminus \{0\}$ .

Denote the random variable  $I_i = (-1)^{\langle \tau, V_i \rangle}$ . Notice that all  $I_1, \dots, I_m$  take values only in  $\{-1, +1\}$ , and the expected value of each  $I_i$  is 0.

Thus, by Hoeffding's inequality there exists a constant  $c \geq 1$  such that:

$$\Pr_S \left[ \left| \frac{1}{m} \sum_{i=1}^m (-1)^{\langle \tau, V_i \rangle} \right| > \epsilon \right] < 2^{-c\epsilon^2 m}$$

If there exists a test  $\tau$  which fails our set  $S$ , then it is not small-biased. Hence, by a union bound over all  $2^n - 1$  tests, we get:

$$\begin{aligned} \Pr_S [\text{S is not small biased}] &\leq \Pr_S [\exists \tau \text{ S fails on } \tau] \leq \\ &\leq \sum_{\tau} \Pr_S \left[ \left| \frac{1}{m} \sum_{i=1}^m (-1)^{\langle \tau, V_i \rangle} \right| > \epsilon \right] < 2^n 2^{-c\epsilon^2 m} \end{aligned}$$

By setting  $m = c^{-1} \cdot \frac{n}{\epsilon^2} = O\left(\frac{n}{\epsilon^2}\right)$ , we can conclude:

$$\Pr_S [\text{S is not small biased set}] < 1$$

Thus, the probability that there exists an  $\epsilon$ -biased set of size  $m$  is strictly greater than zero. That means that there exists an  $\epsilon$ -biased set of size  $m = O\left(\frac{n}{\epsilon^2}\right)$ , as desired.  $\square$

## 11.2 Mergers

### 11.2.1 Motivation

Consider a set of  $r$  random variables  $X_1, \dots, X_r \in \{0, 1\}^n$ , each of which is an  $n$ -bit string. Assume at least one of them is uniformly distributed, but it is not known which one exactly. Moreover, the other random variables are considered "heavily dependent", meaning we cannot assume anything about the correlation between them. For instance, it might be that  $X_2 = X_1$ ,  $X_3 = 2X_7$ ,  $X_4 = -X_8$  etc.

Informally speaking, our goal is to "compress" the  $r$  random variables into a single new random variable, also distributed over  $\{0, 1\}^n$ , while preserving as much of the uniformity as possible.

Much of the work in this field was pioneered by Zeev Dvir ??, which the current lecture is based on.

**Proposition 11.2** (Merger attempt). *Given  $r$  random variables  $X_1, \dots, X_r \in \{0, 1\}^n$ , at least one of which is uniform, we would like to construct an algorithm  $Merg : (\{0, 1\}^n)^r \rightarrow \{0, 1\}^n$  such that  $Merg(X_1, \dots, X_r)$  is uniform as well.*

Using terms from the field of information theory, our merger is supposed to output  $n$  bits of entropy, given a set containing at least  $n$  bits of entropy.

*Example 11.3.* If it were known which of the variables is uniform, our merger would simply use it as its output. For example, in case we knew  $X_3$  is uniform, we'd define our merger to be  $Merg(X_1, \dots, X_r) = X_3$ .

The difficulty of constructing a good merger stems from the fact that we don't know which of the input variables is uniform, and in fact we know nothing about them whatsoever.

Unfortunately, our goal as stated above is somewhat unrealistic - assuming nothing about the input, we wish to perfectly preserve its randomness in the output. We can never expect a truly uniformly random output in this settings.

To make the task at hand more feasible, we will weaken the demands:

First, we will allow ourselves to use some extra bits of randomness in the input.

Second, since we cannot expect our merger to be truly uniform, we will want to be as close to it as possible.

## 11.2.2 Close to Uniform

How can we formally define a distribution which is not uniform per se, but rather close to it?

**Definition 11.4** ( $(\rho, \epsilon)$ -random). Let  $X \in \Omega$  be a random variable over a finite sample space  $\Omega$ .  $X$  is called  $(\rho, \epsilon)$ -random if

$$\forall T \subset \Omega \text{ s.t. } |T| < |\Omega|^\rho \quad \Pr[X \in T] < \epsilon$$

Intuitively, the above definition says that if we try to "capture" a  $(\rho, \epsilon)$ -random variable in a  $\rho$ -fraction of our sample space, we would succeed with very little probability. This means that our random variable is pretty evenly distributed across all the sample space, making it close to uniform.

In fact, there exists an alternative definition for a  $(\rho, \epsilon)$ -random variable, which further emphasizes the resemblance to being uniformly distributed:

Notice that for a uniform random variable  $X \in \Omega$ , it holds that

$$\forall T \subset \Omega \quad \Pr[X \in T] = \frac{|T|}{|\Omega|}$$

Using this characteristic, we can derive

**Definition 11.5** ( $(\rho, \epsilon)$ -random #2). Let  $X \in \Omega$  be a random variable over a finite sample space  $\Omega$ .  $X$  is called  $(\rho, \epsilon)$ -random if

$$\forall T \subset \Omega \text{ s.t. } |T| < |\Omega|^\rho \quad \left| \Pr[X \in T] - \frac{|T|}{|\Omega|} \right| < \epsilon$$

*Observation 11.6.* The two definitions of  $(\rho, \epsilon)$ -randomness are equivalent. Indeed, notice that while in the first definition  $\Pr[X \in T]$  is intended to be close to 0, in the second definition it is close to  $\frac{|T|}{|\Omega|} = |\Omega|^{\rho-1}$ . Since  $|\Omega|^{\rho-1} \rightarrow 0$  as  $\rho \rightarrow 1$ , by adjusting  $\epsilon$  accordingly we can freely transition between the two definitions.

*Remark.* The true definition of  $(\rho, \epsilon)$ -randomness is based on smooth min-entropy from information theory, and will not be presented in the lecture. For our intents and purposes, the definition given above is equivalent and more practical for construction of such  $(\rho, \epsilon)$ -random variables.

### 11.2.3 Definition

Now, after clarifying what it means to be "close to uniform", we can finally properly establish what a Merger is

**Definition 11.7** (Merger). A function  $Merg : (\{0, 1\}^n)^r \times \{0, 1\}^d \rightarrow \{0, 1\}^n$  is called an  $(\rho, \epsilon)$ -merger if given a set of random variables  $(X_1, \dots, X_r)$  distributed over  $\{0, 1\}^n$ , at least one of which is uniform, the random variable  $Merg(X_1, \dots, X_r, Y)$  is  $(\rho, \epsilon)$ -random. The random variable  $Y \sim U_d$  is uniform over  $d$  bits and independent of all the  $X_i$ 's, denoting a pure random seed.

For our algorithm to be considered a "good" merger, we wish to optimize the following

- The seed length  $d$  should be as small as possible
- $\rho \rightarrow 1$ , so the output distribution is close to uniform
- $\epsilon \rightarrow 0$ , so the output distribution is close to uniform

## 11.3 Constructing a Merger

In the next part of the lecture, we will construct a randomness merger using polynomials over finite fields. We will start with the private case of merging two variables, namely  $r = 2$  using the previous notation. Then, we will use the ideas developed in this toy case and generalize them to merge an arbitrary number of variables.

**Theorem 11.8** (Existence of Merger). *For any  $\rho, \epsilon$ , there exists a  $(\rho, \epsilon)$ -random merger which uses a seed of length  $d = O(\frac{1}{\alpha} \log \frac{nr}{\epsilon})$  random bits, where  $\alpha = 1 - \rho$ .*

*Proof.* The proof goes by construction, which we will shortly see in details. □

While the term  $\log \frac{nr}{\epsilon}$  does not bother us so much, the term  $\frac{1}{\alpha}$  is very problematic: Recall that a good merger aims for  $\rho \rightarrow 1$ . So, the better our merger is, more random bits must be invested, as  $\frac{1}{\alpha}$  grows rapidly.

*Note.* Our merger operates on  $(\{0, 1\}^n)^r$ , meaning it receives as input  $r$   $n$ -bit strings. As such, it needs at least  $\log r$  bits to index an input string, and at least  $\log n$  bits to access said string. Thus, we cannot expect to use less than  $O(\max\{\log n, \log r\}) = O(\log n + \log r) = O(\log(nr))$  bits in the seed.

### 11.3.1 Merging 2 random variables

Let  $X, Y \in \{0, 1\}^n$  be two random variables distributed over  $n$ -bit strings, at least one of them is uniformly distributed. Once again, we assume nothing about the correlation between the two variables. Additionally, let  $S \in \{0, 1\}^d$  be our uniform independent random seed over  $d$  bits. Given  $\rho, \epsilon$ , our goal is to build a merger  $Merg(X, Y, S)$  which is  $(\rho, \epsilon)$ -random, using as small as possible seed length  $d$ .

Before we begin our construction, let's play a bit with the input. Instead of looking at  $X, Y \in \{0, 1\}^n$  as  $n$ -bit strings, we can split them up into blocks of size  $b$ , which will be decided at a later point in the proof. There are  $l = n/b$  such blocks.

A block of  $b$  bits is an element of  $\mathbb{F}_2^b$ , which as we know is isomorphic to  $\mathbb{F}_{2^b}$ .

Denote  $q = 2^b$ . Under this new perspective, we can look at  $X, Y$  as vectors in  $(\mathbb{F}_q)^l$ .

As for our seed, it will be beneficial to split it into 2 parts. Meaning, we sample two elements in the field  $\mathbb{F}_q$ , namely  $A, B \in \mathbb{F}_q$ . Our seed  $S$  would now be denoted as  $S = (A, B) \in (\mathbb{F}_q)^2$ , and notice that seed-length =  $2b = 2 \log q$ .

Using the above notation, we define our merger as

$$Merg(X, Y, S) = AX + BY = (AX_1 + BY_1, \dots, AX_l + BY_l)$$

Need to show that  $Merg$  is  $(\rho, \epsilon)$ -random. Specifically

$$\forall T \subset (\mathbb{F}_q)^l \text{ s.t. } |T| < |(\mathbb{F}_q)^l|^\rho = q^{(1-\alpha)l} \quad \Pr_{A,B,X,Y} [AX + BY \in T] < \epsilon$$

Let  $T \subset (\mathbb{F}_q)^l$  be such a group.

Trying to work with  $T$  proves to be quite hard, since it can be pretty much anything.

We take a different approach then, using the **polynomial method**:

Instead of reasoning about the set  $T$ , we will cover it with a bigger and easier to work with group  $T'$ . Since  $T \subset T'$  it holds that

$$\Pr[AX + BY \in T] \leq \Pr[AX + BY \in T']$$

so it is enough to show that  $\Pr_{A,B,X,Y}[AX + BY \in T'] < \epsilon$ . The name of this method comes from the fact that this new group  $T'$  will be defined as the set of roots of a polynomial - an entity which we are very fond of in this course.

Before continuing further, we take a short detour in our proof, to create a powerful

tool which will help us a lot in the parts to come.

### 11.3.1.1 Schwartz-Zippel lemma

*Remark.* Recall the total degree of a polynomial is the maximal of the sums of all the powers of the variables in one single monomial.

For example:  $\deg(x_1^2x_2x_3^4 + 4x_2^5 - 2x_1x_3) = \max\{7, 5, 2\} = 7$

**Theorem 11.9** (Schwartz-Zippel). *Let  $f \in \mathbb{F}[x_1, \dots, x_n]$  be a non-zero polynomial of total degree  $d \geq 0$  over a field  $\mathbb{F}$ . Let  $S$  be a finite subset of  $\mathbb{F}$  and let  $r_1, r_2, \dots, r_n$  be selected at random independently and uniformly from  $S$ . Then*

$$\Pr[f(r_1, r_2, \dots, r_n) = 0] \leq \frac{d}{|S|}$$

In its essence, the Schwartz-Zippel theorem is a generalization of the fundamental theorem of algebra for multivariable polynomials.

*Proof.* The proof goes by induction on  $n$ . The base case where  $n = 1$  is derived directly from the fundamental theorem of algebra, by which we know the polynomial  $f(x)$  has at most  $\deg(f)$  roots. This gives us the base case. Now, assume that the theorem holds for all polynomials in  $n - 1$  variables. We can consider  $f$  to be a polynomial in " $x_1$  only" by rewriting it as

$$f(x_1, \dots, x_n) = \sum_{i=0}^d x_1^i f_i(x_2, \dots, x_n)$$

Since  $f$  is non-zero, there is some  $i$  such that  $f_i$  is also non-zero. We take the largest such  $i$ . Notice that  $\deg(f_i) \leq d - i$ , since  $\deg(x_1^i f_i) \leq d$ .

Now, we randomly pick  $r_2, \dots, r_n$  from  $S$ . By the induction hypothesis

$$\Pr[f_i(r_2, \dots, r_n) = 0] \leq \frac{d - i}{|S|}$$

Note that if  $f_i(r_2, \dots, r_n) \neq 0$ , the univariate polynomial  $f(x_1, r_2, \dots, r_n)$  is of degree  $i$  (and in particular, non-zero) so

$$\Pr[f(r_1, r_2, \dots, r_n) = 0 | f_i(r_2, \dots, r_n) \neq 0] \leq \frac{i}{|S|}$$

Denote the event  $f(r_1, r_2, \dots, r_n) = 0$  by  $A$ , and the event  $f_i(r_2, \dots, r_n) = 0$  by  $B$ . Using conditional probabilities on the event  $B$ , we have

$$\Pr[A] = \Pr[B] \Pr[A|B] + \Pr[B^C] \Pr[A|B^C] \leq^* \Pr[B] + \Pr[A|B^C] = \frac{d-i}{|S|} + \frac{i}{|S|} = \frac{d}{|S|}$$

where in (\*) we used the fact that any probability is at most 1. □

**Corollary 11.10** (Schwartz-Zippel for finite fields). *Let  $f \in \mathbb{F}_q[x_1, \dots, x_n]$  be a non-zero polynomial over a finite field  $\mathbb{F}_q$ . Thus,  $f$  has at most  $q^{n-1} \deg(f)$  roots in the field.*

*Proof Sketch.* Take  $S = \mathbb{F}_q$ . From Schwartz-Zippel, we know that if  $r_1, \dots, r_n$  are selected independently and uniformly in the field, then

$$\Pr[f(r_1, \dots, r_n) = 0] = \frac{\#roots}{q^n} \leq \frac{\deg(f)}{q} \implies \#roots \leq q^{n-1} \deg(f)$$

□

### 11.3.1.2 The polynomial method

Armed with [Corollary 11.10](#), let's return to showing  $\Pr_{A,B,X,Y}[AX + BY \in T] < \epsilon$  using the polynomial method. Formulating our intentions, we would like a polynomial  $Q \in \mathbb{F}_q[z_1, \dots, z_l]$  with the following properties:

1.  $Q \neq 0$
2.  $\forall t \in T \ Q(t) = 0$  (abbreviated as  $Q|_T = 0$ )
3.  $Q$  has a low degree

Constructing such polynomial can be trivially achieved by interpolating the elements of  $T$ . However, such technique would result in a high degree, violating the third property. Therefore, we will show the existence of such polynomial in a different (and much more interesting) way, using a combinatorial approach.

**Lemma 11.11.** *Given a non-negative integer  $d$  such that:*

$$\binom{d+l}{l} > |T|$$

there exists a polynomial  $Q \in \mathbb{F}_q[z_1, \dots, z_l]$  of degree  $d$  with the aforementioned three properties.

*Remark.* Notice that  $d$  gets larger as  $|T|$  grows (as one would naturally expect).

*Proof.* Given some  $t = (t_1, \dots, t_l) \in T$ , requiring  $Q(t) = 0$  translates to solving a linear constraint on the coefficients of  $Q$ . Moreover, the number of possible monomials in  $Q$  is at most  $\binom{d+l}{l}$  (as the monomial is of degree  $d$  in  $l$  variables). Since the total number of linear constraints is  $|T|$ , which is strictly smaller than the number of unknowns, there is a nontrivial solution, yielding our desired polynomial  $Q$ .  $\square$

### 11.3.1.3 *Merg* is $(\rho, \epsilon)$ -random

As  $Q|_T \equiv 0$ , we can cover our original arbitrary set  $T$  with a smoother one  $T'$ , defined as follows:

$$T \subset T' = \{(z_1, z_2, \dots, z_l) \in (\mathbb{F}_q)^l \mid Q(z_1, \dots, z_l) = 0\}$$

so it's enough to show that  $\Pr_{A,B,X,Y}[AX + BY \in T'] < \epsilon$ .

In our settings, at least one of  $X, Y$  is uniformly distributed, so assume w.l.o.g  $X \sim U_n$  is uniform over  $n$  bits. Denote a sample  $x \sim X$  as *bad* if the following holds:

$$\Pr[AX + BY \in T \mid X = x] \geq \frac{\epsilon}{2}$$

In other words, a sample  $x \sim X$  is *bad* if it contributes a lot to our target probability. Denote the set of all *bad* samples as  $B(X)$ . We shall now prove a simple (yet meaningful) lemma:

**Lemma 11.12.**  $\Pr_{A,B,X,Y}[AX + BY \in T] \geq \epsilon \implies \Pr_{x \sim X}[x \in B(X)] \geq \frac{\epsilon}{2}$

*Proof.* Using the Law of Total Probability, we can write:

$$\begin{aligned} \epsilon &\leq \Pr_{A,B,X,Y}[AX + BY \in T] = \\ &\Pr[AX + BY \in T \mid x \in B(X)] \cdot \Pr_{x \sim X}[x \in B(X)] + \\ &\Pr[AX + BY \in T \mid x \notin B(X)] \cdot \Pr_{x \sim X}[x \notin B(X)] \leq \\ &\leq^* \frac{\epsilon}{2} + \Pr_{x \sim X}[x \in B(X)] \end{aligned}$$

Where  $(*)$  can be simply derived from the following three (trivial) statements:



1.  $\Pr[Ax + BY \in T \mid x \in B(X)] \leq 1$
2.  $\Pr[Ax + BY \in T \mid x \notin B(X)] \leq \frac{\epsilon}{2}$
3.  $\Pr_{x \sim X}[x \notin B(X)] \leq 1$

□

Now, assume by contradiction that our merger fails on  $T$ , meaning

$$\Pr_{A,B,X,Y}[AX + BY \in T] \geq \epsilon$$

Using the previous lemma, we can sample a *bad*  $x \sim X$ , and a  $y$  from the support of  $Y \mid X = x$ , such that:

$$\Pr_{A,B}[Ax + By \in T] \geq \frac{\epsilon}{2}$$

That is, we sample a  $y$  which exploits the "badness" of  $x$ . Notice that now the randomness comes only from the seed, as  $x, y$  are fixed.

Due to the fact that  $T \subset T'$ , we may deduce:

$$\Pr_{A,B}[Ax + By \in T'] \geq \frac{\epsilon}{2}$$

Recall that  $Q$  is a polynomial in  $l$  variables, of total degree  $d$ . If we plug in  $Ax + By = (Ax_1 + By_1, \dots, Ax_l + By_l)$  (where  $x$  and  $y$  are fixed),  $Q$  becomes a polynomial on  $A, B \in \mathbb{F}_q$ .

Denote it by  $R_{x,y}(A, B) = Q(Ax + By)$  of total degree

$$\deg(R_{x,y}) \leq \deg(Q) \leq d$$

From the Schwartz-Zippel lemma, if  $R_{x,y} \not\equiv 0$ ,

$$\#\text{Roots of } R_{x,y} \leq d \cdot q$$

And due to our choice of  $x, y$ ,

$$\Pr_{A,B}[R_{x,y}(A, B) = 0] \geq \frac{\epsilon}{2}$$

Thus,

$$\frac{\epsilon}{2} \cdot q^2 \leq \#\text{Roots of } R_{x,y} \leq d \cdot q$$

We may pick  $q > \frac{2d}{\epsilon}$  so that the above inequality cannot hold.  
Hence  $R_{x,y} \equiv 0$  in  $\mathbb{F}_q[A, B]$ .

As  $Q(x) = R_{x,y}(1, 0) \equiv 0$ , we get that every *bad* sample of  $x \sim X$  is a root of  $Q$ . By applying the Schwartz-Zippel lemma again and using [Lemma 11.12](#),

$$\frac{\epsilon}{2} \cdot q^l \leq \#\text{Roots of } Q \leq d \cdot q^{l-1}$$

which implies  $q \leq \frac{2d}{\epsilon}$ , in contradiction to how we chose  $q$ .

Hence,  $Q \equiv 0$ , in contradiction to our construction of  $Q$  as a non-zero polynomial.  $\square$

### 11.3.1.4 Fixing Parameters

Along the way, we gathered two requirements:

1.  $\binom{d+l}{l} > |T| = q^{(1-\alpha)l}$
2.  $q > \frac{2d}{\epsilon}$

Using the well-known fact:

$$\binom{d+l}{l} \geq \left(\frac{d}{l}\right)^l$$

We shall enforce:

$$\left(\frac{d}{l}\right)^l > q^{(1-\alpha)l} \Rightarrow d > l \cdot q^{(1-\alpha)}$$

Fix  $d = 2l \cdot q^{(1-\alpha)}$ . The second demand implies

$$q > \frac{2d}{\epsilon} = \frac{4l}{\epsilon} \cdot q^{(1-\alpha)}$$

That is,

$$q^\alpha > \frac{4l}{\epsilon} \implies q > \left(\frac{4l}{\epsilon}\right)^{\frac{1}{\alpha}}$$

So we can pick  $q$  to be this threshold. Ultimately, our seed length can be bounded:

$$d = 2 \cdot \log q = O\left(\frac{1}{\alpha} \log \frac{n}{\epsilon}\right)$$

## 11.3.2 Merging multiple random variables

We would like to generalize the last construction to the general case, where there are more than 2 random variables. Assume now there are  $k$  random variables  $X_1, \dots, X_k$  in need of merging, at least one of which is uniform as before.

### 11.3.2.1 Iterative pair merging

One trivial way to construct a merger using the results of the private case of 2 variables is the following construction:

Using an iterative process, at each stage  $i$  the merger receives as input a seed  $S_i$ , and merges each pair of random variables sequentially. The number of random variables is reduced by half after each iteration, so there are  $\log k$  iterations. Although we cannot guarantee there exists a uniform random variable after the first merging iteration, we can ensure there is a  $(\rho, \epsilon)$ -random one, which is sufficient.

It is easy to see that our algorithm would use a total seed-length of  $O(\log k \lceil \frac{1}{\alpha} \log \frac{n}{\epsilon} \rceil)$ .

### 11.3.2.2 Curve Merger

In the two variables scenario, our merger was based on the plane spanned by  $X$  and  $Y$ , namely  $AX + BY$ . This "hyperplane merger" would perform poorly in the case of  $k$  variables, since it would require  $k$  seeds.

Another possibility is to look at the line between  $X$  and  $Y$ , so our merger would now be of form  $(1 - A)X + AY$ . Notice that a point on the line can be indexed by only one parameter, which makes it generalizable: given  $k$  variables, our merger would output points on the curve that passes through them all. Since any point on the curve can be indexed by a single parameter, our seed length would not grow as much as before with regards to  $k$ .

So, let's start by building a curve that passes through each input point  $x_1, \dots, x_k \in \mathbb{F}_{2^n}$ . Choose some  $\gamma_1, \dots, \gamma_k \in \mathbb{F}_{2^n}$  distinct field elements. For each element let us build the following polynomial:

$$c_i(u) = \prod_{j \in [k], j \neq i} \frac{(u - \gamma_j)}{\gamma_i - \gamma_j}$$

we can see that  $c_i(\gamma_j) = 0$  if  $i \neq j$  and 1 if  $i = j$ .

Look at the curve  $C(y) = \sum_{i=1}^k c_i(y) \cdot x_i$ . we can see that each point  $x_1, \dots, x_k$  is on it since plugging  $\gamma_j$  will result in  $C(\gamma_j) = x_j$ . Moreover, since all  $c_i$ 's are of degree  $k - 1$ , so is the curve.

Our merger is defined by the curve:

$$M(X_1, \dots, X_k, Y) = \sum_{i=1}^k c_i(Y) \cdot X_i$$

This merger is called a "curve merger" since it uses the seed to choose a point uniformly on the curve passing through the input.

It is based on the work of Dvir & Wigderson ?.

**Theorem 11.13** (Curve Merger). *For every  $\rho = 1 - \alpha > 0$  the curve merger is  $(\rho, \epsilon)$ -random with a seed of size  $d = O(\log k + \log n)$  and  $\epsilon = O(\frac{1}{nk})$ .*

*Proof.* Pick  $q = 2^d$  such that  $(nk)^{\frac{4}{\alpha}} < q \leq 2(nk)^{\frac{4}{\alpha}}$ .

We will assume w.l.o.g that  $n = d \cdot r$  (losing a small number of bits is negligible).

As before, we can view each  $X_i$  as an element of  $(\mathbb{F}_q)^r$  and the seed as an element of  $\mathbb{F}_q$ . Let  $\epsilon = 4 \cdot q^{-\frac{\alpha}{4}}$ .

Note that  $\epsilon = O(\frac{1}{nk})$  and  $d = \log q = O(\log n + \log k)$ .

Assume by contradiction that  $Z = M(X_1, \dots, X_k, Y)$  fails and is not  $(\rho, \epsilon)$ -random. Explicitly speaking, there exists some  $T \subset (\mathbb{F}_q)^r$  s.t.  $|T| \leq q^{r\rho} = 2^{n(1-\alpha)}$  for which  $\Pr[Z \in T] \geq \epsilon$ .

Denote  $s = q^{1-\frac{\alpha}{2}}$ . Observe that, since  $r < n < q^{\frac{\alpha}{4}}$ , we have

$$\binom{s}{r} \geq \left( \frac{q^{1-\frac{\alpha}{2}}}{q^{\frac{\alpha}{4}}} \right)^r \geq q^{r(1-\alpha)} \geq |T|$$

Next, we utilize the polynomial method to analyze our merger. As before, we will build a non-zero polynomial  $Q$  whose group of roots includes the whole set  $T$ . To prove such a polynomial exists, we resort to linear algebra:

The number of monomials of degree  $s$  in  $r$  variables is

$$\binom{s}{r} = \binom{s+r-1}{r} \geq \binom{s}{r} \geq |T|$$

Looking at the requirement  $Q|_T = 0$  as a system of linear equations, we have more degrees of freedom than constraints. Therefore we can build a non-zero polynomial  $Q \in \mathbb{F}_q[z_1, \dots, z_r]$  with degree  $\leq s$  such that  $Q(t) = 0$  for every  $t \in T$ , simply by solving the system of linear equations.

Our new goal is to show that  $Q$  has too many roots, which is a contradiction to its low degree.

Assume w.l.o.g that  $X_1$  is the uniform variable (the proof will be identical if another source is uniform).

We say that  $x_1$  is "bad" iff  $\Pr[Z \in T | X_1 = x_1] \geq \frac{\epsilon}{2}$ .

As seen before, by an averaging argument  $\Pr[X_1 \text{ is bad}] \geq \frac{\epsilon}{2}$ .

**Lemma 11.14.** *If  $x_1$  is bad then  $Q(x_1) = 0$ .*

*Proof.* Fix some bad  $x_1$ .

Since  $\Pr[Z \in T | X_1 = x_1] \geq \frac{\epsilon}{2}$ , there are  $x_2, \dots, x_k$  such that

$$\Pr[Z \in T | X_1 = x_1, \dots, X_k = x_k] \geq \frac{\epsilon}{2}$$

So the randomness comes only from the seed.

Observe the curve defined by the  $x_i$ 's

$$C = \left\{ \sum_{i=1}^k c_i(u) \cdot x_i \mid u \in \mathbb{F}_q \right\} = \{M(x_1, \dots, x_k, u) \mid u \in \mathbb{F}_q\}$$

The restriction of  $Q$  to the curve is given by the univariate polynomial:

$$R(u) = Q \left( \sum_{i=1}^k c_i(u) \cdot x_i \right)$$

By definition of  $c_i$ 's the curve's degree is  $k - 1$ , so by composition  $\deg(R) \leq s \cdot (k - 1)$ .

We saw that  $\Pr[Z \in T | X_1 = x_1, \dots, X_k = x_k] \geq \frac{\epsilon}{2}$ , and since  $T$  is contained in the roots of  $Q$ , we deduce that  $R$  is zero on at least  $\frac{\epsilon}{2}$ -fraction of  $\mathbb{F}_q$ .

Using the inequality:

$$\frac{s \cdot (k - 1)}{q} < \frac{s \cdot (q^{\frac{\alpha}{4}} - 1)}{q} < \frac{\epsilon}{2}$$

$R$  is zero on at least  $\frac{\epsilon}{2}q > s \cdot (k - 1)$  points.

By the Schwartz-Zippel lemma,  $R$  has at most  $\deg(R) \leq s \cdot (k - 1)$  roots, which is a contradiction - unless  $R$  is the zero polynomial. Hence, assume that  $R$  is indeed the zero polynomial, so it is zero on all of  $\mathbb{F}_q$ .

In particular  $\gamma_1$  is also a root, which means  $0 = R(\gamma_1) = Q(x_1)$  as required.  $\square$

To finish the proof, lets get another contradiction using Schwartz-Zippel:

We've already proved that  $\Pr[X_1 \text{ is bad}] \geq \frac{\epsilon}{2}$  therefore at least  $\frac{\epsilon}{2}$ -fraction of  $(\mathbb{F}_q)^r$  are bad, which means that the number of roots  $Q$  has is at least

$$\frac{\epsilon}{2} \cdot |\mathbb{F}_q^r| > \frac{s}{q} \cdot |\mathbb{F}_q^r| = \frac{s}{q} \cdot q^r = s \cdot q^{(r-1)}$$

By Schwartz-Zippel  $Q$  has at most  $\deg(Q) \cdot q^{r-1} \leq s \cdot q^{r-1}$  roots, which is a contradiction!

Therefore, our initial assumption is disproved, meaning that the curve merger  $M(X_1, \dots, X_k, Y)$  is  $(\rho, \epsilon)$ -random.  $\square$

## LECTURE 12

# TWO SOURCE EXTRACTORS AND UNBALANCED EXPANDERS

---

In this lecture, we'll talk about expander graphs. Expander constructions have spawned research in pure and applied mathematics, with several applications to complexity theory, design of robust computer networks, and the theory of error-correcting codes. Before we dive into the details, we'll take a small detour and introduce a different topic: Two-Source Extractors.

## 12.1 Two Source Extractors

We would like to explore the possibility to create a method that, given two independent non-uniform random stream-sources, "produces" a purely uniform bit of data. Such method will be called "Two-Source Extractor". This is highly useful for many fields which require an unbiased random source but only have biased (non-uniform) sources. In order to define this more formally, we require the notion of min entropy which we will define in the next section.

### 12.1.1 Definitions

**Definition 12.1** (min-entropy  $k$ ). Let  $X$  be a random variable over the sample space  $\{0, 1\}^n$ .  $X$  is said to have **min-entropy**  $k$  if  $X$  is uniform over some (unknown)  $S \subseteq \{0, 1\}^n$ , such that  $|S| = 2^k$

**Definition 12.2** (Two-Source Extractor). A two source extractor for min-entropy  $k$  is a function  $\text{Ext}: \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$  s.t. for any **independent**  $X, Y$  random variables with min-entropy  $k$  (each),  $\text{Ext}(X, Y)$  has a small bias, for example smaller than 0.01.

*Remark.* We note that we can require the bias to be as small as  $O\left(2^{-\frac{k}{2}}\right)$  but we simplify for our purposes.

We can also look at the definition above from a graphical point of view. We can view  $\text{Ext}$  as a matrix  $A \in \mathbb{F}_2^{2^n \times 2^n}$  such that  $A_{ij} = \text{Ext}(i, j)$ . In this setting, the requirement

above is translated to the fact that for any subset of cells of size bigger or equal to  $2^k \times 2^k$  (we note that the cells does not have to be contiguous) the bias of the values of the cells is small. The conditions can be illustrated by the illustration bellow.

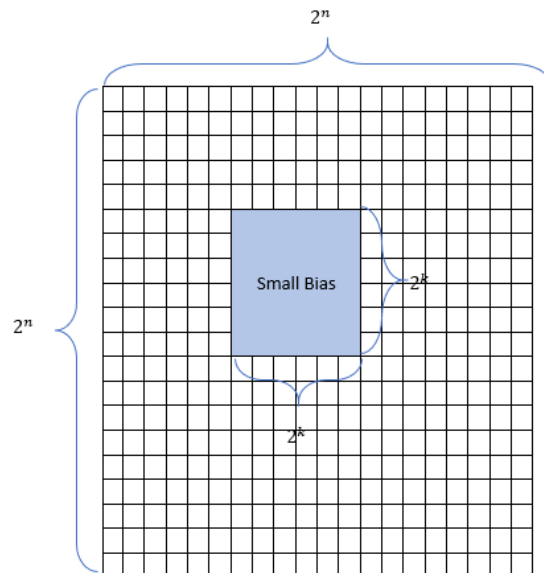


Figure 8: An illustration of a two-source extractor. A two-source extractor can be viewed as an  $2^n \times 2^n$  truth table. Our demand implies that no area of sufficient size can not have large bias (the area doesn't have to be contiguous).

The idea was first proposed by Benny Shor and Oded Regev. It was shown that solving this problem could solve an open problem in graphs, proposed by Ramsey. Among other uses for these extractors is to extract randomness from sources that were leaked but still have some entropy in them.

## 12.1.2 Main Results And Constructions

We first note that it is possible to prove (though we won't show it) that for any  $n \in \mathbb{N}$  there exist such extractor for  $k = \log_2 n + O(1)$ . Alas, this proof is non constructive. We note that the hardness of the problem is since the extractor only get a sample for each of the random variables, and doesn't have any knowledge of the underlying distribution. We also note that any randomly sampled matrix of size  $2^n \times 2^n$  will be



a two-source extractor, but such construction is not efficient in any way.

We will present several constructions which were achieved during the years:

- $k = \frac{n}{2} + O(1)$  - In their original article Shor and Reg propose  $Ext(X, Y) = \langle X, Y \rangle = \sum_i x_i y_i \pmod 2$  where  $X, Y \in \mathbb{F}_2^n$ .
- $k = 0.499999n$  - In Bourgain [2005] it was shown the first explicit two-source extractor with ratio of less than 0.5. The construction can be described by  $Ext(X, Y) = \langle X^3 + X, Y^3 + Y \rangle$  where  $X, Y \in \mathbb{F}_{2^n}$ . The proof yield from a fact we mentioned,  $\max(|A + A|, |A \cdot A|) \geq |A|^{1+\epsilon}$  but the proof is completely out of our scope. The proof was a breakthrough and inspired many other results.
- State of the art - The best result known to this day is for  $k = (\log n)^{\log \log \log n}$ .
- Paley Matrix - We look at finite field  $\mathbb{F}_p$  and  $X, Y$  in it. Then, if  $X + Y$  is a square of some element in  $\mathbb{F}_p$  then  $Ext(X, Y) = 1$  and vice versa. This construction is believed to uphold with  $k = O(\log n)$  but no proof has been presented.

## 12.2 Unbalanced Expanders

### 12.2.1 Introduction

Let us now return to the subject of unbalanced expanders. Let us define the notion of an unbalanced expander graph.

**Definition 12.3** (Unbalanced Expander). Given  $n, k, \epsilon$  an unbalanced expander is a bipartite graph  $G = (L, R, E)$  such that:

1.  $G$  is left  $d$ -regular
2.  $|L| = n$
3.  $\forall S \subseteq L$  such that  $|S| \leq k$  it holds that  $|\Gamma(S)| \geq (1 - \epsilon) \cdot d \cdot |S|$ .

*Remark.* We note that our goal will be to minimize  $d$  and  $|R|$  (as a function of  $n, k, \epsilon$ ). Also, we will denote  $|R|$  by  $m$  henceforth.

Bellow, an illustration of the definition is presented. We note that our demand is that even for small sets  $S$  the neighbor set  $\Gamma(S)$  is not too small.

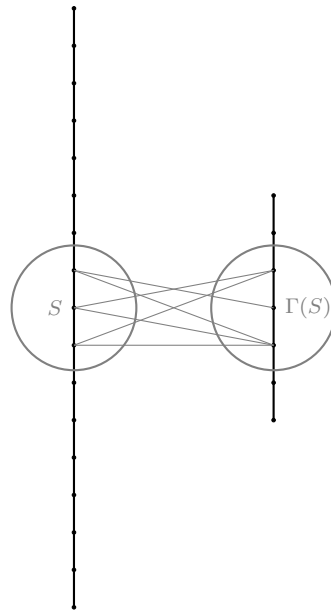


Figure 9: An illustration of the definition of Unbalanced Expander. We wish to ensure that  $\Gamma(S)$  is not too small

The motivation for this construction is to compress entropy. Consider a very large  $n = |L|$  and a very small  $m = |R|$ . Then our construction takes a random variable that is distributed only on  $|S|$ , and creates a new random variable that is distributed on  $\Gamma(S)$  (i.e. the neighbors of  $S$  in  $G$ ). The first random variable has little entropy relative to  $n$ , the size of its probability space, as  $|S| \ll n$ . But the new random variable has a large entropy relative to  $m$ , the size of its probability space.

## 12.2.2 Proof Of Existence

We will show a proof of existence of unbalanced expanders via the probabilistic method, with "optimal"  $d, m$  (thus, the proof is non constructive).

**Theorem 12.4.**  $\forall n, k, \varepsilon$  there exists an unbalanced expander with:

$$d = \Theta \left( \log \left( \frac{n}{k} \right) \right)$$

$$m = \Theta\left(c^{\frac{1}{\varepsilon}} dk\right)$$

*Proof.* Fix  $S \subseteq L$ ,  $|S| = k$  and  $T \subseteq R$ ,  $|T| \leq (1 - \varepsilon)dk - 1$ . We'll bound  $\Pr[\Gamma(S) \subseteq T]$  and then we'll use union bound over all  $S, T$  to conclude that the total probability  $< 1$ . Thus there exist a graph which is a valid unbalanced expander.

In our random left  $d$ -regular graph, each edge is drawn uniformly and independently. So we have  $dk$  edges from  $S$  and every edge has probability of  $\frac{|T|}{m}$  to fall in  $T$  (we ignore repeated edges as this only reduces the probability and we want to bound it). Hence,

$$\Pr[\Gamma(S) \subseteq T] \leq \left(\frac{|T|}{m}\right)^{dk}$$

Taking the sum over all  $S, T$  by a union bound we have an upper bound of

$$\binom{n}{k} \binom{m}{(1 - \varepsilon)dk} \left(\frac{(1 - \varepsilon)dk}{m}\right)^{dk} = (*)$$

as we have  $\binom{n}{k}$  options for  $S$  and  $\binom{m}{(1 - \varepsilon)dk}$  options for  $T$ . If  $(*) < 1$ , we're done. We'll use the inequality

$$\binom{n}{k} \leq \left(\frac{ne}{k}\right)^k \leq \left(\frac{3n}{k}\right)^k$$

and get that,

$$(*) \leq \left(\frac{3n}{k}\right)^k \left(\frac{3m}{(1 - \varepsilon)dk}\right)^{(1 - \varepsilon)dk} \left(\frac{(1 - \varepsilon)dk}{m}\right)^{dk}$$

To complete the proof, we would like to show that

$$\left(\frac{3n}{k}\right)^k \left(\frac{3m}{(1 - \varepsilon)dk}\right)^{(1 - \varepsilon)dk} < \left(\frac{m}{(1 - \varepsilon)dk}\right)^{dk}$$

We note that it suffices to remove the  $(1 - \varepsilon)$  in  $3^{(1 - \varepsilon)d}$  as it only helps the bound, and take the  $k$  root of both sides,

$$3^d \frac{3n}{k} \left(\frac{m}{(1 - \varepsilon)dk}\right)^{(1 - \varepsilon)d} < \left(\frac{m}{(1 - \varepsilon)dk}\right)^d$$

Also, it suffices to ignore  $(1 - \varepsilon)$  in the denominator,

$$3^d \frac{3n}{k} < \left(\frac{m}{dk}\right)^{\varepsilon d}$$

Taking  $m \triangleq 6^{\frac{1}{\varepsilon}} dk$ , we get

$$3^d \frac{3n}{k} < 6^d$$

$$\Downarrow$$

$$\frac{3n}{k} < 2^d$$

and thus

$$d = \Theta \left( \log \left( \frac{n}{k} \right) \right)$$

□

*Remark.* We note that the proof only bound the probability for groups of size exactly  $k$ , and we need to bound it for groups of size up to  $k$ . To complete the proof we need to sum over  $|S| \leq k$ , as considered above only  $|S| = k$ . To that end, we can show with the same parameters that  $(*) < 4^{-|S|}$  for a specific  $|S|$  (since we took the  $k$ th root of the inequality, this will only add a constant factor). Summing over all  $|S| \leq k$  would get us to  $(*) < 1$  as we required

*Remark.* If we don't wish to compress the entropy, i.e. we don't require that  $m$  will be smaller than  $n$ , we can achieve that with a constant  $d$ , where  $d = \Theta \left( \log \left( \frac{n}{m} \right) \right)$

### 12.2.3 An Explicit Polynomial Based Construction

We will show the following theorem.

**Theorem 12.5.** *For any  $\alpha > 0$ ,  $N, K, \varepsilon > 0$  there exists an explicit construction for an unbalanced expander such that  $D = \Theta \left( \left( \frac{\log N \log K}{\varepsilon} \right)^{1+\alpha} \right)$  and  $M = \Theta(D^2 K^{1+\alpha})$*

*Remark.* We first note that we will show the main ideas of the construction given in [Guruswami et al. \[2006\]](#). We therefore will be a bit informal and will not show how to derive the parameters of the construction from the given parameters (in particular we will not fully define  $q, h$  which will be described bellow). For full details one can refer to [Guruswami et al. \[2006\]](#)

We look at a field of size  $q = D, \mathbb{F}_q$ . We will derive below a condition for the size of  $q$ . Each vertex  $v \in L$  is associated with a polynomial  $f_v$  with degree smaller than  $n$  (we denote this by  $f_v \in \mathbb{F}_q^{<n}[x]$ ). Thus our left side is of size  $N = q^n$ . Our right side will consist of vectors of size  $m+1$  over  $\mathbb{F}_q$ . Hence, the right side is of size  $M = q^{m+1}$ .

For every vertex  $v \in L$  we associate each edge with an element  $x \in \mathbb{F}_q$  such that

$$e_x^v = \left( v, \left( x, f_v(x), f_v^h(x), \dots, f_v^{h^{m-1}}(x) \right) \right)$$

Where  $h^m = K$ . This implies that our graph is left  $D = q$  regular. We note that other constraints will be derived during the construction.

It is important to clarify what does it mean to exponent a polynomial in our context. Each  $f_v$  is a polynomial of degree less than  $n$ . We note that this is a vector space over  $\mathbb{F}_q$ . As we've seen before this space is isomorphic (as vector spaces) to the field  $\mathbb{F}_q[x]/\langle p(x) \rangle$  where  $p(x)$  is some irreducible polynomial of degree  $n$ . Since  $\mathbb{F}_q[x]/\langle p(x) \rangle$  is a field the multiplication is well defined in it. Thus, the definition of our multiplication is to multiply the two isomorphic elements in the field and to use this element as the result (this is similar to what was done in lecture 11).

We note here that it is enough to show that for any set  $T \subseteq R$  such that  $|T| \leq (1 - \varepsilon)DK - 1$  then  $|\Gamma^{-1}(T)| \leq K - 1$  where  $\Gamma^{-1}(T) = \{v \in L \mid \Gamma(v) \subseteq T\}$ . Let there be a set  $T$  of size  $|T| \leq (1 - \varepsilon)DK - 1$ . We use the polynomial method to construct a polynomial "wrapping"  $T$  which will be easier to work with. We search for a polynomial  $Q$  such that:

1.  $Q$  is not the zero polynomial
2.  $Q(t) = 0$  for any  $t \in T$
3.  $Q$  has a "low degree" - we require that  $\deg_{x_0} Q \leq (1 - \varepsilon)q = (1 - \varepsilon)D$  (where  $\deg_{x_0}$  is the degree of polynomial in  $x_0$  alone) and  $\deg_{x_i} Q = h - 1$  for  $i = 1, \dots, m - 1$
4.  $Q$  is of minimal degree in  $x_0$  of the polynomial which hold the previous items

We note that the fourth item might seem a bit arbitrary at the moment but will be crucial for the proof. We want to show that such  $Q$  does exist.

**Lemma 12.6.** *There exist such a  $Q$  with the above properties*

*Proof.* We note that if we show a non empty set of  $Q$ s that hold items 1-3, choosing the minimal such one is always possible. Thus, we need to show that there exists a

non empty set of  $Q$  that holds items 1-3. We look at all of the polynomials of the degree described in item 3. Since each  $x_0$  is up to power  $(1 - \varepsilon)D$  and each other  $x_i$  is up to power  $h^m$  we note that the dimension of that vector space is

$$(1 - \varepsilon) Dh^m = (1 - \varepsilon)DK$$

Each point in  $T$  adds a linear constraint on this space, that reduces the dimension of the space by one. Thus, the space of polynomials that hold the items 2 and 3 is

$$(1 - \varepsilon)DK - T > 1$$

Thus, the space is not the zero space (its dimension is bigger than 1). Hence, there exist a non empty set of  $Q$ s that holds item 1-3, as wanted.  $\square$

We look at a vertex  $v$  such that  $\Gamma(v) \subseteq T$  (meaning  $v$  is a "bad" vertex, which is constrained in  $T$ ). This vertex is associated with  $f_v$  a polynomial. Hence, for any  $x \in \mathbb{F}_q$  it holds that

$$\left( x, f_v(x), f_v^h(x), \dots, f_v^{h^{m-1}}(x) \right) \in T$$

Thus, by our choice of  $Q$  it holds that,

$$Q \left( \left( x, f_v(x), f_v^h(x), \dots, f_v^{h^{m-1}}(x) \right) \right) = 0$$

We note that we can look at  $Q \left( \left( x, f_v(x), f_v^h(x), \dots, f_v^{h^{m-1}}(x) \right) \right)$  as polynomial in  $x$ . We denote

$$R_{f_v}(x) = Q \left( \left( x, f_v(x), f_v^h(x), \dots, f_v^{h^{m-1}}(x) \right) \right)$$

Thus, for any  $x \in \mathbb{F}_q$   $R_{f_v}(x) = 0$  for a "bad"  $v$ . Alas, it holds that,

$$\deg R_{f_v} \leq \deg_{x_0} Q + \sum_i \deg_{x_i} Q \cdot \deg f_v \leq (1 - \varepsilon)q - 1 + (h - 1)(n - 1)m <$$

$$(1 - \varepsilon)q - 1 + hnm \stackrel{(1)}{<} (1 - \varepsilon)q - 1 + q\varepsilon < q$$

Where in (1) we note that it is enough to require that  $q > \frac{hmn}{\varepsilon}$  for it to hold. Hence, since  $R_{f_v}$  has  $q$  roots (the elements of the field) and degree smaller than  $q$  it holds that  $R_{f_v}(x) = 0$

We now return to the fact that each of our vertices is associated with a polynomial of degree smaller than  $n$ . As we stated above, these polynomials are elements in a field  $\mathbb{F}_q[x]/\langle p(x) \rangle$ . With this in mind we can think of  $Q$  as  $Q^*$  such that,

$$Q^*(z) \in \mathbb{F}_q[x]/\langle p(x) \rangle$$

$$Q^*(z) = Q(x, z, z^h, z^{h^2}, \dots, z^{h^{m-1}})$$

Were  $x$  is just the element of field that the natural projection assigns to  $x$ . With such a view, we can see that  $f_v$  for a "bad"  $v$  is a root of  $Q^*$  since,  $Q^*(f_v) = R_{f_v} = 0$ . Hence if  $Q^*$  **isn't the zero polynomial** it can have at most roots as his degree, which will give us a bound of number of "bad" vertices. We note that the degree of  $Q^*$  holds,

$$\deg_z Q^* \leq \sum_{i=1}^m \deg_{x_i} Q \cdot \deg_z z^{h^{i-1}} = (h-1) \cdot (1 + h + h^2 + \dots + h^{m-1}) =$$

$$h^m - 1 = K - 1$$

Hence, since we showed that  $Q^*$  has at most  $K - 1$  many roots, and each "bad" vertex is a root it holds that,  $|\Gamma^{-1}(T)| \leq K - 1$  if  $Q^*$  is not the zero polynomial. Thus all that is left to show is that  $Q^*$  is **not** the zero polynomial.

**Lemma 12.7.**  *$Q^*$  is not the zero polynomial*

*Proof.* We look at  $Q$

$$Q(x_0, x_1, \dots, x_m) = \sum_{I=(i_1, i_2, \dots, i_m)} g_I(x_0) x_1^{i_1} \cdot \dots \cdot x_m^{i_m}$$

Assume in contradiction that for all  $I$   $g_I \in \langle p(x) \rangle$  (which implies  $Q^*$  is the zero polynomial). Thus, we can look at

$$\bar{Q}(x_0, \dots, x_m) = \sum_{I=(i_1, i_2, \dots, i_m)} \frac{g_I(x_0)}{p(x_0)} x_1^{i_1} \cdot \dots \cdot x_m^{i_m}$$

It is clear that  $\deg_{x_0} \bar{Q} < \deg_{x_0} Q$  and also that  $\bar{Q}(T) = 0$  since  $p$  is irreducible (meaning it has no roots in  $\mathbb{F}_q$ ) and  $\bar{Q} = \frac{1}{p(x_0)} Q(x_0, \dots, x_m)$ . Also, it is clear the  $\bar{Q}$  is indeed a valid polynomial since  $g_I \in \langle p(x) \rangle$  for any  $I$ . Since  $\bar{Q}$  is a contradiction to how we chose  $Q$  (it upholds all the conditions and has a smaller degree in  $x_0$ ), there

exists such an  $I$  such that  $f_I \notin \langle p(x) \rangle$ . We now note that for every  $I \neq I'$

$$x_1^{i_1} \cdots x_m^{i_m} \neq x_1^{i'_1} \cdots x_m^{i'_m} \quad \text{for } x_i = z^{h^{i-1}}$$

Since we count in  $h$  base, and each representation in the base is unique. Hence no other element can cancel the  $I$  such that  $I \notin \langle p(x) \rangle$  and thus,  $Q^*$  is not the zero polynomial.  $\square$

Thus, we have shown the main ideas for which the construction is based on. From this point it remains to show that the parameters for the construction (i.e  $n, q, h, m$ ) from the given parameters (i.e  $N, K, \varepsilon, \alpha$ ) and that a irreducible polynomial can be constructed efficiently. These steps are shown in [Guruswami et al. \[2006\]](#)



## References

1. J. Bourgain. More on the sum-product phenomenon in prime fields and its applications. *International Journal of Number Theory*, 1(1), 2005.
2. Guruswami, Umans, and Vadhan. Extractors and condensers from univariate polynomials. In *Technical Report TR06-134, Electronic Colloquium on Computational Complexity*, 2006.
3. Ian Nicholas Stewart. *Galois theory*. CRC Press, 2015.