# MASTERARBEIT

Titel der Masterarbeit

## Additive generation of rings by units with particular reference to Dedekind domains

verfasst von

## Thomas Blank BSc

angestrebter akademischer Grad

## Master of Science (MSc)

# Contents

**Introduction**

## 1.1 Introduction I and preliminaries

The subject of decomposing a ring element, specifically a rational integer, into a sum of elements featuring a particular property is ancient. One of the earliest studies concerning such decompositions are related to Pythagorean triples. Later on there arose the question of writing any integer as a sum of two squares - the famous solution being due to Fermat. Whilst it does not take much mathematical ability to formulate these problems, the expertise needed to tackle them is tremendous, especially when the summands are required to have properties related to the notion of primality.

Rather than asking for decompositions of elements into sums of primes or squares this thesis deals with the topic of additively representing a ring element as sum of invertible elements of the ring. Clearly, the issue is trivial for rational integers, but appears to be everything but obvious in other rings.

Plainly because the notion of a ring originated much later than the use of rational integers in mathematical history, one is not mistaken by suspecting the discipline of additive unit representations to be young; indeed it goes back to the 1950s, when Zelinsky [1] proved that every element of an endomorphism ring of a vector space over a division ring is expressible as sum of two automorphisms. Later, Henriksen [2] established that no matter the ring, any matrix is 3-good; the terminology $k$-good being due to Vámos [3] meaning representable as sum of *exactly* $k$ units of the matrix ring. The problem of a ring satisfying that all its elements are additively expressible as a sum of units of a given length, led to the notion of the unit sum number: One defines the *unit sum number* $u(S)$ of a ring $S$ to be

$$u(S) = \begin{cases} k & S \text{ is } k\text{-good, but not } j\text{-good for all } j < k \text{ with } j, k \in \mathbb{N} \\ \omega & S \text{ is not } k\text{-good for any } k \in \mathbb{N}, \text{ but every element is a finite sum of units} \\ \infty & \text{there exists an element in } S \text{ not expressible as a finite sum of units in S} \end{cases},$$

where we say the ring $S$ is $k$-good, if every element of $S$ is $k$-good.

Investigations about the unit sum number problem for matrix rings quickly led to the application of matrix normal forms for commutative rings and an interest in certain types of rings affiliated to these normal forms such as elementary divisor rings and Hermite rings. Owing to this method, it was proved that the ring of square matrices over a principal ideal domain possesses unit sum number two, a result stemming from the Smith normal form for matrix rings over principal ideal domains.

In 1972, Levy [4] was able to obtain an astonishing relation between the unit sum number of a matrix ring over a Dedekind domain $\mathfrak{O}$ and its class number $h_{\mathfrak{O}}$, exhibiting the first application of the unit sum number results for matrix rings developed earlier. For his proof he relied heavily on the work of Steinitz [5] about the structure of finitely generated modules over Dedekind domains, which is proved by utilising the theory of projective modules. The matrix-theoretic part is based on the work of Krull [6] about block matrices.

Still related to matrix rings, but of a very different flavour, is the unit sum number problem in the case of non-commutative rings. The issue is completely settled in the case of Artinian rings or, more generally, semilocal rings. The theorem classifying the behaviour of these rings with respect to their unit sum number is composed of Zelinsky's result mentioned above and the renowned Artin-Wedderburn structure theorem for semisimple rings.

The most natural rings to consider are perhaps the rings of algebraic integers of number fields, since they constitute the direct generalisation of rational integers. The case of quadratic number fields was already fully dealt with by Belcher in 1974, later certain cubic and quartic fields were examined. However, only classes of number fields featuring a single fundamental unit have been successfully explored thus far. Though no general method to determine the unit sum number of a ring of algebraic integers is known, a major break-through was achieved by Jarden and Narkiewicz [7] in 2005. They proved that a finitely generated integral domain can only have unit sum number $\omega$ or $\infty$. The techniques used are a deep result about the number of solutions to unit equations by Evertse, Schlickewei and Schmidt [8] based on Schmidt's subspace theorem and the classical van der Waerden's theorem [9] about arithmetic progressions. Refining the methods by employing Szemerédi's theorem [10], Jarden and Narkiewicz demonstrated that the density of

$$N_n = \{x \in \mathbb{N} | x \text{ is } k\text{-good in } K \text{ for some } k \leq n\},$$

where $K$ denotes a number field, is zero in $\mathbb{N}$.

## 1.2 Notation and conventions

### Rings

As our objective is the study of ring elements being expressible as sums of units, we require every occurring ring to be *associative* and *unital*, i.e. featuring a neutral multiplicative element. Moreover we require ring homomorphisms $S \to S'$ to map $1_S \mapsto 1_{S'}$.

As we will encounter commutative and non-commutative rings featuring profoundly differing theories, we use the symbol $R$ *for commutative rings* and $S$ *whenever the ring is also allowed to be non-commutative*. A ring free of zero-divisors is called a *domain*, a commutative domain $\mathcal{D}$ will also be referred to as *integral domain*. In this paper all occurring *principal ideal domains* are understood to be commutative. Moreover we write $\mathfrak{O}$ or $\mathcal{O}$, whenever we want to indicate that a ring is a *Dedekind domain* or a *ring of algebraic integers* respectively.

Given a ring $S$, we denote by $S^*$ its *group of units*, by $\mathrm{Jac}(S)$ we understand the *Jacobson radical* of $S$, i.e. the intersection of all maximal left ideals of S.[1]

Other than explicitly mentioned we require every ideal to be non-zero. When working in Dedekind rings we want the notion of *ideals* also to incorporate *fractional ideals* - if we refer to ideals in the ordinary sense, we will speak of *integral ideals*. Moreover the *ideal class* of a non-zero ideal $M$ will be denoted by $[M] = M + P_K$, where $P_K$ is the subgroup of principal fractional ideals within the group of all fractional ideals of some Dedekind domain.

The *opposite ring* of a ring $(S, +, \cdot)$ is given by $(S^{\mathrm{op}}, +, \circ)$, where $\forall a, b \in S^{\mathrm{op}} : a \circ b = b \cdot a$.

Eventually note that other than in $\mathbb{Z}$ the notion of a *greatest common divisor* $\gcd()$ is only defined up to associated elements, however, this will be negligible in our calculations.

### Matrices

By the *size* of a matrix we mean the number of its columns and rows $m \times n$; we will also frequently say a matrix has size $n$, if it is a square $n \times n$ matrix. By $\mathrm{Mat}_n(S)$ we denote the matrix ring of square matrices of size $n$ over $S$.

As usual $I_n$ denotes the identity matrix of size $n$. We will frequently solely write $I$ for the identity matrix and $\mathbf{0}$ for a matrix featuring only zero entries, if the size is apparent from the context.

---

[1]Note that substituting right for left ideals in the definition, does not alter the resulting ideal. See [11, Theorem 13.8].

A matrix $A$ with entries from a commutative ring is defined to be *singular* if $\det(A) = 0$, else we say $A$ is *regular*. The notion of a matrix $A$ being *unimodular*, being *contained in* $\mathrm{GL}_n(S)$ or being a *unit* in $\mathrm{Mat}_n(S)$ all signify $\det(A) \in S^*$, whence $A$ admits an inverse in $\mathrm{Mat}_n(S)$.

A *nilpotent matrix* $A$ satisfies the existence of a positive integer $k$, such that $A^k$ equals the zero-matrix.

### Number fields

A *number field* is a field $K \subseteq \mathbb{C}$ admitting a finite $\mathbb{Q}$-basis. By $[K : \mathbb{Q}]$ we denote the degree of $K$ as a $\mathbb{Q}$-vector space. The *ring of algebraic integers* affiliated to $K$ will be denoted by $\mathcal{O}_K$. We use $\mathrm{Tr}_K()$ and $\mathcal{N}_K()$ to symbolise the *field trace and norm* respectively. The *discriminant* of an element $x$ will be symbolised by $d_K(x)$, whereas $d_K$ signifies the *field discriminant* of $K$.

## 1.3 Introduction II

We are interested in the question, whether a ring has the property that each element is representable as a finite sum of units. Adopting the terminology of Ashraf and Vámos [3] we call an element $x$ of some ring $S$ $k$-*good*, if

$$\exists (\eta_i)_{i=1}^k \in (S^*)^k : x = \sum_{i=1}^{k} \eta_i.$$

We will not use the notion 1-good, as those elements simply correspond to units. A subset of a ring is called $k$-good, if all elements within this subset are $k$-good. Furthermore we define the *unit sum number* $u(S)$ to be

$$u(S) = \begin{cases} k & S \text{ is } k\text{-good, but not } j\text{-good for all } j < k \text{ with } j, k \in \mathbb{N} \\ \omega & S \text{ is not } k\text{-good for any } k \in \mathbb{N}, \text{ but every element is a finite sum of units} \\ \infty & \text{there exists an element in } S \text{ not expressible as a finite sum of units in S} \end{cases} .$$

It is easily shown that the subset $S^\omega := \{\sum_{i=1}^k \epsilon_i | k \in \mathbb{N}, \epsilon_i \in S^*\}$ is a subring of $S$. In particular $u(S) = \infty \Leftrightarrow S^\omega \subsetneqq S$. Moreover $u(S) > 1$ for any unital ring S, as $0 \in S \setminus S^*$.

Our first objects of examination are fields. The next lemma implicates that for all fields $F$, we have $F = F^\omega$.

**Lemma 1.1.** *Let $F$ be a field then $u(F) = 2$, unless $F$ is isomorphic to the field of two elements $\mathbb{F}_2$. In this case $u(F) = u(\mathbb{F}_2) = \omega$.*

*Proof.* If $F \not\cong \mathbb{F}_2$, $F$ must contain at least three elements, thus satisfying

$$\forall x \in F : \exists y \neq x \in F^* : x = (x - y) + y \text{ and } x - y \in F^*.$$

For $F \cong \mathbb{F}_2$ it is evident that 0 and 1 can not be represented by unit sums of equal length $k \in \mathbb{N}$ and thus $u(\mathbb{F}_2) = \omega$. $\qquad \square$

**Remark 1.2.** Note that a ring homomorphism maps units to units: Suppose $\epsilon \in S^*$ and $\varphi(\epsilon) = a \notin S'^*$, then $1_{S'} = \varphi(\epsilon)\varphi(\epsilon^{-1})$. Hence $\varphi(\epsilon^{-1})$ is the inverse of $a$ in $S'$.

We collect a few general, simple facts about $k$-good rings, that will be useful later.

**Lemma 1.3** ([12, Lemma 1,2])**.**

(i) *A non-trivial ring epimorphism maintains the ring-properties of being $k$-good or $\omega$-good.*

(ii) *Let $I$ be an ideal of a ring $S$ contained in $Jac(S)$. Then $x \in S$ is $k$-good, if under the canonical epimorphism its image in $S/I$ is $k$-good. The converse also holds.*

*(iii) Let $\{S_i\}_{i=1}^r$ be a finite family of rings, where each $S_i$ constitutes a $k_i$-good ring. Then $\prod_{i=1}^r S_i$ is $k$-good, where $k = \max\{k_i | 1 \leq i \leq r\}$.*

*Proof.*

(i) Let $\varphi : S \to S'$ denote a ring epimorphism and $b \in S'$. There exists $a \in S$ such that $\varphi(a) = b$, which admits a representation $a = \sum_{i=1}^k \eta_i$, where $\eta_i \in S^*$. Then $\varphi(a) = \sum_{i=1}^k \varphi(\eta_i)$, which is a sum of $k$ units in $S'$.

(ii) Let $\varphi : S \to S/I$ denote the canonical epimorphism. For some $x \in S$ let $\overline{x} := \varphi(x)$ be $k$-good, i.e. $\overline{x} = \sum_{i=1}^k \overline{e_i}$, where $\overline{e_i} \in S/I^*$. Since there exist $\overline{d_i} \in S/I^*$, such that $e_i d_i + I = 1 + I$, this entails $1 - e_i d_i \in I \subseteq \mathrm{Jac}(S)$. Hence

$$1 + (-1)(1 - e_i d_i) \in S^*,$$

which shows that the $e_i$'s are in fact units in $S$.
As $y := x - \sum_{i=1}^k e_i \in \mathrm{Jac}(S)$, we see that $1 + e_1^{-1} y \in S^*$ and whence $e_1 + y \in S^*$. This implies that

$$x = (e_1 + y) + \sum_{i=2}^k e_i$$

is $k$-good in $S$.

(iii) Set $k = \max\{k_i | 1 \leq i \leq r\}$. Take an arbitrary $r$-tuple $(x_1, \ldots, x_r) \in \prod_{i=1}^r S_i$. Then $x_i' := x_i - (k - k_i)1$ lies in $S_i$ and is therefore $k_i$-good. Thus $x_i = x_i' + \sum_{j=k_i+1}^k 1$ shows that any element $x_i$ in the $i$-th component of $\prod_{i=1}^r S_i$ is also $k$-good. Now as all $x_i$ are $k$-good, we finish with

$$(x_1, \ldots, x_r) = \Big( \sum_{j=1}^k \epsilon_{1j}, \ldots, \sum_{j=1}^k \epsilon_{rj} \Big) = \sum_{j=1}^k (\epsilon_{1j}, \ldots, \epsilon_{rj})$$

for suitable $\epsilon_{ij} \in S_i^*$ - the last term being a sum of $k$ units in $\prod_{i=1}^r S_i$.

$\square$

**Example 1.4.** By listing all rings of order four[2], it can be proved that $\mathbb{F}_2 \oplus \mathbb{F}_2 := \mathbb{F}_2 \oplus \mathbb{F}_2$ with componentwise multiplication is the smallest ring possessing an infinite unit sum number.

**Example 1.5.** Since $\mathbb{Z}^* = \{\pm 1\}$, it is evident that $u(\mathbb{Z}) = \omega$. (cf. 6.15)

---

[2]See [13].

### 1.3.1 Extension of the length of an additive unit representation

Before attending to specific rings, we discuss the possibilities of extending the length by which an element is expressible as sums of units.

Raphael [14] calls a ring *even*, if its neutral multiplicative element 1 can be written as sum of an even number of units - otherwise *odd*. The definition is equivalent to demanding the zero element to be representable as sum of an odd number of units. Note that the smallest odd ring is the field of two elements. We also introduce a slightly finer distinction namely $k$-even to indicate the length $k$ needed to express 1. It is useful to note that for a ring $S$ we have for even $k$: $S^*$ is $k$-good $\Leftrightarrow$ $S$ is $k$-even.

**Lemma 1.6** (cf. [14])**.**

(i) *The ring-property $(k$-$)$even is invariant under non-trivial ring homomorphisms; particularly for rings $S_1, \ldots, S_n$ we have $\prod_{i=1}^n S_i$ is even, iff every $S_i$ is even.*

(ii) *Let the ring $S_1$ be $k$-even. Let $S_2$ be odd or satisfy that $1_{S_2}$ is $j$-good for some $j > k$, but not $\ell$-good for some $\ell \leq k$. Then $\mathrm{Hom}(S_1, S_2) = \{0\}$.*

(iii) *Let $S_i$ be rings satisfying $u(S_i) = \omega$, at most one of them not even. Then $u(\prod_{i=1}^k S_i) = \omega$.*

*Proof.*

(i) Trivial.

(ii) Immediate from (i).

(iii) W.l.o.g. let $S_1$ be odd. We prove the assertion for $k = 2$, as the argument is easily extendable to the general case by induction. Both zero elements $0_{S_1}$ and $0_{S_2}$ are representable as sums of two units, since $0 = 1 - 1$. As $S_2$ is even, there exists an odd number $\kappa$, such that the zero element of $S_2$ may be written as sum of $\kappa$ units. Now taking an arbitrary $(a, b) \in S_1 \times S_2$ we know that $a = \sum_{i=1}^{k_1} \epsilon_i, b = \sum_{i=1}^{k_2} \eta_i$ for some $\epsilon_i \in S_1^*, \eta_i \in S_2^*$ and $k_1, k_2 \in \mathbb{N}_{>1}$. Consider the two progressions

$$P_1(s) = 2s + k_1 \text{ and } P_2(s, t) = 2s + \kappa t + k_2.$$

As $P_1 \cap P_2$ is not empty, we find $K$ being the minimal positive integer in $P_1 \cap P_2$. The proof is complete by writing

$$(a, b) = \Big( \sum_{i=1}^{k_1} \epsilon_i, \sum_{i=1}^{k_2} \eta_i \Big) = \sum_{i=1}^{K} (\epsilon_i, \eta_i)$$

for certain units $\{\epsilon_i\}_{i=k_1+1}^K \subseteq S_1, \{\eta_i\}_{i=k_2+1}^K \subseteq S_2$, that fulfil

$$\sum_{i=k_1+1}^K \epsilon_i = 0_{S_1} \text{ and } \sum_{i=k_2+1}^K \eta_i = 0_{S_2}.$$

$\square$

**Example 1.7.** Let $\mathcal{D} \neq \mathbb{F}_2$ be an integral domain satisfying that $1_\mathcal{D}$ is not 2-good and let $K$ be its field of fractions. Then (ii) of the previous lemma assures $\text{Hom}(K, \mathcal{D}) = \{0\}$, since $K$ clearly is 2-even (cf. Lemma 1.1).

**Remark 1.8.**

(i) As seen the beneficial property of an even ring $S$ is, that sums of units may be "filled up" to increase their length; more precisely an element $a = \sum_{i=1}^k \epsilon_i \in S$, $\epsilon_i \in S^*$, is $n$-good for any $n \geq k$.

(ii) If the ring is not assumed to be even, there is in general no need for a $k$-good element to be $j$-good for some specific $j \neq k$: As an example take $1 \in \mathbb{F}_2$. Clearly 1 is 3-good, but neither 2-good or 4-good.

The behaviour of whole rings compared to single elements differs, when it comes to extensions of the sum length.

**Lemma 1.9.** *Let $S$ be a $k$-good ring, then $S$ is also $n$-good for all $n \geq k$. In particular $S$ is either $k$-even or $(k+1)$-even.*

*Proof.* Take an arbitrary $x \in S$. There exist units $\epsilon_1, \ldots, \epsilon_k$ in $S$, such that $x = \sum_{i=1}^k \epsilon_i$. Now $h := \sum_{i=1}^{k-1} \epsilon_i$ is contained in $S$ as well and is therefore $k$-good. Thus we see that $x = h + \epsilon_k$ is $(k+1)$-good. Induction verifies the statement for all $n \geq k$.

The second part follows easily: If $k$ is even, then we can write 1 as sum of $k$ units, thus $S$ is $k$-even. If $k$ is odd, 1 is also $(k+1)$-good, due to the first part. Hence $S$ is $(k+1)$-even.

$\square$

The property *even* in the context of rings of algebraic integers will be discussed in Chapter 6.

# The unit sum number of matrix rings

Henriksen [2] in 1972 exhibited that a matrix ring $\mathrm{Mat}_n(S)$ over an arbitrary ring $S$ fulfils $u\big(\mathrm{Mat}_n(S)\big) \in \{2,3\}$. This result led to new interest in the unit sum number problem for matrix rings and paved the way for further developments and techniques, such as the use of specific types of rings and the application of matrix normal forms.

Many of the upcoming proofs demand the usage of block (diagonal) matrices:

**Definition 2.1** ([15, Definition 1.1])**.** Let $A_1$ be an $m_1 \times n_1$ matrix and $A_2$ be an $m_2 \times n_2$ matrix. By the *block diagonal sum* $\mathrm{diag}(A_1, A_2)$ we mean the so called *block diagonal matrix*

$$\begin{pmatrix} A_1 & \mathbf{0} \\ \mathbf{0} & A_2 \end{pmatrix} \text{ of size } (m_1 + m_2) \times (n_1 + n_2).$$

It is useful to extend this definition to "zero-size" matrices, to avoid the consideration of many special cases later on: For any $m \times n$ matrix $(a_{ij})$ we set

$$\mathrm{diag}\big((a_{ij}), \mathbf{0}_{(0 \times p)}\big) = \begin{pmatrix} a_{11} & \cdots & a_{1n} & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} & 0 & \cdots & 0 \end{pmatrix},$$

symbolising $(a_{ij})$ with $p$ columns appended. By $\mathrm{diag}\big((a_{ij}), \mathbf{0}_{(p \times 0)}\big)$ we refer to the same construct appending $p$ rows instead of columns.

It is convenient to list some basic computational rules.

**Remark 2.2.**

(i) Let $A, B$ be matrices split up into blocks $A_{ij}, B_{ij}$:

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \text{ and } B = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}.$$

Suppose all the blocks are compatible in shape to perform the standard row-by-column matrix multiplication, then one verifies $C = AB$, where $C_{ik} = \sum_{j=1}^{2} A_{ij} B_{jk}$. We calculate the first block $C_{11}$ of the new matrix to clarify, what is meant by compatible in shape: As $C_{11} = A_{11}B_{11} + A_{12}B_{21}$, we must have the following sizes of the blocks.

| Block | Size |
|-------|------|
| $A_{11}$ | $i \times j$ |
| $B_{11}$ | $j \times k$ |
| $A_{12}$ | $i \times \ell$ |
| $B_{21}$ | $\ell \times k$ |

Evidently this rule of calculation is extendable to matrices decomposed into more than four blocks.

(ii) Let $A = \mathrm{diag}(A_1, \dots, A_r)$ be a block diagonal sum of square matrices $A_i$ over a commutative ring, then one calculates with ease $\det(A) = \prod_{i=1}^{r} \det(A_i)$.

(iii) Note that we are using diag() in two ways. On the one hand we want it to symbolise the block diagonal sum $\mathrm{diag}(A_1, A_2)$. On the other hand, if we are using a single argument $\mathrm{diag}(A)$, we refer to the tuple of entries in the diagonal of $A$.

We remark that (i) signifies that as long as the blocks of the matrices $(A_{ij})$ and $(B_{ij})$ are compatible the multiplication equals the standard matrix multiplication with the elements being matrices themselves.

**Definition 2.3.** We define two matrices $A, B$ to be *equivalent*, denoting $A \sim B$, if there exist invertible matrices $P, Q$ such that $A = PBQ$. Square matrices $A$ and $B$ are called similar, if $A = PBP^{-1}$ for some invertible matrix $P$.

**Definition 2.4.** As usual by *elementary matrix operations* on a matrix $A$ over a ring $S$ we mean either

- *a switching of rows (columns).* For example

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} c & d \\ a & b \end{pmatrix}$$

switches the first and last row of $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$. Evidently the used elementary matrix is invertible, as it is self-inverse.

- *a multiplication of one row (column) by a unit element.* For instance

$$\begin{pmatrix} \epsilon & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} \epsilon a & \epsilon b \\ c & d \end{pmatrix}$$

multiplies the first row of $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$ by a unit $\epsilon$ of $S$. The elementary matrix is invertible, as $\begin{pmatrix} \epsilon^{-1} & 0 \\ 0 & 1 \end{pmatrix}$ is its inverse.

- *addition of a multiple of a row (column) to another.* For example

$$\begin{pmatrix} 1 & r \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} a + rc & b + rd \\ c & d \end{pmatrix}$$

adds an $r$-multiple, $r \in S$, of the last row to the first. The inverse of this elementary matrix is given by $\begin{pmatrix} 1 & -r \\ 0 & 1 \end{pmatrix}$.

These operations correspond to left (right) multiplication with unimodular matrices. In particular applying them to $A$ produces a matrix $A'$, which lies in the same equivalence class $[A]_\sim$ as $A$ does.

## 2.1 Every matrix is $3$-good

Let $S$ denote an arbitrary, unital, not necessarily commutative ring.

**Proposition 2.5** ([2, Lemma 1,2])**.**

 (i) *Diagonal matrices and nilpotent matrices of size $n > 1$ over $S$ are 2-good.*

 (ii) *Every $A \in \mathrm{Mat}_n(S)$ is the sum of a diagonal matrix and an invertible matrix.*

*Proof.*

 (i) Let $D = \mathrm{diag}(a_1, \ldots, a_n)$, then $D = U_1 + U_2$, where

$$U_1 = \begin{pmatrix} a_1 & 0 & 0 & \ldots & 0 & 1 \\ 1 & a_2 & 0 & \ldots & 0 & 0 \\ 0 & 1 & a_3 & \ldots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \ldots & 1 & a_{n-1} & 0 \\ 0 & 0 & 0 & \ldots & 1 & 0 \end{pmatrix}, U_2 = \begin{pmatrix} 0 & 0 & 0 & \ldots & 0 & -1 \\ -1 & 0 & 0 & \ldots & 0 & 0 \\ 0 & -1 & 0 & \ldots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \ldots & -1 & 0 & 0 \\ 0 & 0 & 0 & \ldots & -1 & a_n \end{pmatrix}.$$

We show by elementary matrix operations, that $U_1$ is equivalent to the identity matrix $I$; First clear out the entry $a_1$ in position $(1, 1)$ by adding a suitable multiple of the last column. Now successively clear out entries $a_2, a_3, \ldots$ by adding a suitable multiple of the first, second, etc. column. The matrix, we

obtain, is of the form

$$\begin{pmatrix} 0 & 0 & \ldots & 0 & 1 \\ 1 & 0 & \ldots & 0 & 0 \\ 0 & 1 & \ldots & 0 & 0 \\ \vdots & \vdots & \ddots & 0 & 0 \\ 0 & 0 & \ldots & 1 & 0 \end{pmatrix}.$$

Clearly multiplication with an appropriate permutation matrix shows that $U_1 \sim I$ and hence $U_1$ is invertible. The procedure to show that $U_2 \sim I$ is of similar simplicity, yielding that $U_2$ is invertible.

For a nilpotent matrix $N$ it suffices to note that $(I - N)^{-1} = \sum_{i \geq 0} N^i$, where only finitely many summands do not vanish, is the inverse of $I - N$ and hence invertible. Thus $N = (I - N) + I$ yields a decomposition into two units.

(ii) For the second part note that the assertion holds trivially for $n = 1$ as $a = (a - 1) + 1$ for any $a \in S$. Assume the hypothesis holds for matrices of fixed size $n$. Any $A' \in \mathrm{Mat}_{n+1}(S)$ is of the form $A' = \begin{pmatrix} A & \mathbf{b} \\ \mathbf{c} & \delta \end{pmatrix}$ with $A \in \mathrm{Mat}_n(S)$, $\delta \in S$ and $\mathbf{b}, \mathbf{c}$ vectors of length $n$ over $S$. Now $A = D + U$, $D$ being a diagonal matrix, $U$ an invertible matrix. Define

$$D' = \begin{pmatrix} D & \mathbf{0} \\ \mathbf{0} & \delta - 1 - \mathbf{c}U^{-1}\mathbf{b} \end{pmatrix} \text{ and } U' = \begin{pmatrix} U & \mathbf{b} \\ \mathbf{c} & 1 + \mathbf{c}U^{-1}\mathbf{b} \end{pmatrix}$$

so that $A' = D' + U'$. Let $I \in \mathrm{Mat}_n(S)$ be the identity matrix. We need to prove that $U'$ is invertible; define

$$P = \begin{pmatrix} I & \mathbf{0} \\ -\mathbf{c}U^{-1} & 1 \end{pmatrix} \text{ and } Q = \begin{pmatrix} U^{-1} & -U^{-1}\mathbf{b} \\ \mathbf{0} & 1 \end{pmatrix}$$

now $PU'Q$ equals the identity matrix of $\mathrm{Mat}_{n+1}(S)$ due to Remark 2.2(i). The proof is complete by showing that $P, Q \in GL_{n+1}(S)$:
For $P$ use the last column to clear out all entries $-\mathbf{c}U^{-1}$, which entails $P \sim I$.
For $Q$ use the last row to clear the entries $-U^{-1}\mathbf{b}$ to get $\begin{pmatrix} U^{-1} & 0 \\ 0 & 1 \end{pmatrix}$, its inverse being $\begin{pmatrix} U & 0 \\ 0 & 1 \end{pmatrix}$.

$\square$

The proposition allows us to draw an immediate conclusion:

**Corollary 2.6** ([2, Theorem 3]). *Let $n > 1$, then the matrix ring $\mathrm{Mat}_n(S)$ over any ring $S$ is 3-good.*

For this reason the unit sum number of any matrix ring of size greater one can only take the values two or three. Observe also that the statement of the corollary does *not* imply $u(\text{Mat}_n(S)) = 3$, because $\text{Mat}_n(S)$ being 3-good does not hinder it to be 2-good as well. Our next task is to find classes of rings, which induce 2-good matrix rings.

## 2.2 Types of rings associated with matrix rings

We introduce some classes of rings useful in dealing with the unit sum number problem of matrix rings. For the sake of simplicity, we restrict ourselves to the case of unital, commutative rings $R$.

**Definition 2.7.** As we are not necessarily dealing with quadratic matrices, we define a matrix $(a_{ij})$ of any size to be *lower triangular*, if $a_{ij} = 0$, whenever $i < j$. The definition of *upper triangular* runs similarly. $(a_{ij})$ is *diagonal*, if $a_{ij} = 0$, whenever $i \neq j$.

- A *Bézout ring* $R$ is a ring satisfying, that the sum of two principal ideals is again principal.

- $R$ is called *Hermite ring*, if every matrix of size greater one over $R$ is equivalent to a lower triangular matrix.

- Kaplansky [16] calls a ring *elementary divisor ring*, if every $m \times n$ matrix $A$ of size greater one admits a *diagonal reduction*: there exists an integer $r$ and a matrix $D = \text{diag}(d_1, \ldots, d_r, \mathbf{0}_{(m-r)\times(n-r)})$ with $d_i | d_{i+1}$, $1 \leq i \leq r - 1$, such that $A \sim D$.

**Remark 2.8.** There are numerous results and conditions under which one of the above type of ring equals another, we can merely present a selection of them here; for a thorough synopsis the reader shall be pointed to the work of Lam [17, §4].

For an example of a Bézout ring that is not an elementary divisor ring see [18, § 3 and § 4].

We state the two main facts about the dependencies among the classes of rings just introduced.

**Theorem 2.9** ([17, cf. Theorem 4.25]). *For a commutative ring $R$ we have the following chain of implications.*

$R$ *is a principal ideal domain* $\Rightarrow$ $R$ *is an elementary divisor domain* $\Rightarrow$
$R$ *is a Hermite domain* $\Rightarrow$ $R$ *is a Bézout domain.*

*Proof.* To prove the first implication we will employ the Smith normal form, Theorem 2.22, which will be treated in the next section.

For the second suppose $R$ is an elementary divisor ring, then $R$ being a Hermite ring is immediate from the definition, since every diagonal matrix is also a triangular matrix.

To prove that a Hermite ring is a Bézout ring take $(b_1, b_2) \in R^2$. As $R$ is Hermite

$$\exists Q \in \mathrm{GL}_2(R) \text{ and } \exists d \in R \text{ such that } Q(b_1, b_2)^\mathsf{T} = (0, d)^\mathsf{T}.$$

This shows that $d$ is a linear combination of $b_1, b_2$ and hence $dR \subseteq b_1 R + b_2 R$. To obtain the other inclusion, let $(x_1, x_2)^\mathsf{T}$ be the second column of $Q^{-1}$, now

$$(b_1, b_2)^\mathsf{T} = Q^{-1}(0, d)^\mathsf{T} = d(x_1, x_2)^\mathsf{T} \text{ thus } b_1 R + b_2 R \subseteq dR.$$

$\square$

**Remark 2.10.** It is easy to show that a commutative Noetherian Bézout domain $\mathcal{D}$ is in fact a principal ideal domain: Since $\mathcal{D}$ is Noetherian every ideal $\mathfrak{a}$ is finitely generated; $\mathfrak{a} = (a_1, \ldots, a_n) \subseteq \mathcal{D}$. A repeated application of the defining property, that the sum of two principal ideals is again principal, implicates the principality of $\mathfrak{a}$. Explicitly $(a_1, \ldots, a_n) = \sum_{i=1}^n a_i R = aR = (a)$, for a suitable $a \in \mathcal{D}$.

Next we state a result, which provides a sufficient and necessary condition for a Noetherian integral domain to be a principal ideal domain by means of matrix decompositions.

**Proposition 2.11** ([19, Proposition 3.1 and 3.2]). *A Noetherian integral domain $\mathcal{D}$ is a principal ideal domain, iff for all $A \in \mathrm{Mat}_2(\mathcal{D})$ there is an invertible matrix $U$ and a symmetric matrix $S$ with $A = SU$.*

*Proof.*

"$\Leftarrow$" Let $a, b \in \mathcal{D}$ be arbitrary and set $A = \begin{pmatrix} a & 0 \\ b & 0 \end{pmatrix}$. There exists $U = (u_{ij})_{1 \leq i, j \leq 2} \in \mathrm{GL}_2(\mathcal{D})$, such that $AU$ is symmetric. Now

$$\begin{pmatrix} a & 0 \\ b & 0 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \end{pmatrix} = \begin{pmatrix} au_{11} & au_{12} \\ bu_{11} & bu_{12} \end{pmatrix}$$

yields the condition $au_{12} = bu_{11}$. Denote $\epsilon := \det(U) = u_{11}u_{22} - u_{12}u_{21} \in \mathcal{D}^*$, then

$$a\epsilon = au_{11}u_{22} - au_{12}u_{21} = u_{11}(au_{22} - bu_{21}),$$
$$b\epsilon = bu_{11}u_{22} - bu_{12}u_{21} = u_{12}(au_{22} - bu_{21}).$$

This leads to $(au_{22} - bu_{21}) \subseteq (a, b) \subseteq (au_{22} - bu_{21})$, implying that $(a, b)$ is principal and we have shown that $\mathcal{D}$ is a Bézout domain. Finally Remark 2.10 yields the claim.

"⇒" For $A \in \mathrm{Mat}_2(\mathcal{D})$ we obtain in view of Theorem 2.9 matrices $P, Q \in \mathrm{GL}_2(\mathcal{D})$, such that $PAQ = D$ is diagonal. Setting

$$E = \begin{pmatrix} \det(P) & 0 \\ 0 & \det(Q)^{-1} \end{pmatrix}$$

one sees that $PAQE = DE$ is still diagonal. Whence

$$AQE(P^{-1})^\intercal = P^{-1}(PAQE)(P^{-1})^\intercal = P^{-1}DE(P^{-1})^\intercal$$

is symmetric and by construction $\det(QE(P^{-1})^\intercal) = 1$, thus invertible.     □

**Remark 2.12.** From Proposition 2.5 we conclude, that every matrix over an elementary divisor ring is the sum of two units. A fortiori by Theorem 2.9 also matrix rings over principal ideal domains, in particular euclidean domains, have unit sum number two. Thus apart from attempting to prove the 2-good property directly, we may also prove that certain rings are elementary divisor rings.

## 2.3  Matrix normal forms

In the former section we used matrix equivalence to define Hermite rings and elementary divisor domains. The purpose of normal forms - which exist for both types of rings - is to uniquely link a matrix to a specific triangular or diagonal matrix by equivalence. The reader interested in an extensive treatment of integral matrices, i.e. matrices with coefficients in a principal ideal domain, may find the book of Newman [20] valuable.

**Proposition 2.13** ([16, Theorem 3.5]). *A matrix $A$ over a Hermite ring admits a so called Hermite normal form, which warrants the existence of a unimodular matrix $U$, such that $AU$ is lower triangular.*

*Proof.* At first let $A$ be a $1 \times n$ matrix. For $n = 2$ the definition of Hermite ring grants the existence of invertible $P, Q$ such that $PAQ = (d\ 0)$. If we require $Q$ to be a $2 \times 2$ matrix, $P$ needs to be a scalar, which may thus be omitted. Now suppose the assertion holds for $1 \times (n-1)$ matrices. Let $A$ be a $1 \times n$ matrix over $R$, we find a $1 \times (n-1)$ matrix $B$ such that $A = (a\ B)$. Due to the induction hypothesis there is

$$V \in \mathrm{GL}_{n-1}(R) \text{ satisfying } BV = (b\ 0 \ldots 0).$$

Furthermore there exists $W \in \mathrm{GL}_2(R)$ with $(a\ b)W = (d\ 0)$. Putting

$$U = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & V \end{pmatrix} \begin{pmatrix} W & \mathbf{0} \\ \mathbf{0} & I_{n-2} \end{pmatrix},$$

we see $AU = (d\ 0 \ldots 0)$ by Remark 2.2(1) .

To treat the general case let $A$ be an $m \times n$ matrix. We find an invertible matrix V, which reduces the first row: $AV = \begin{pmatrix} a & \mathbf{0} \\ \mathbf{b} & C \end{pmatrix}$. By induction on $m$ the $(m-1) \times (n-1)$ matrix C gives rise to $W \in \mathrm{GL}_{n-1}(R)$ with $CW$ lower triangular. Therefore

$$AV \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & W \end{pmatrix} = \begin{pmatrix} a & \mathbf{0} \\ \mathbf{b} & C \end{pmatrix} \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & W \end{pmatrix} = \begin{pmatrix} a & \mathbf{0} \\ \mathbf{b} & CW \end{pmatrix},$$

which is lower triangular.                                                      □

**Remark 2.14.**

- Note that the statement may be rephrased as: $A$ is equivalent to a lower triangular matrix, the left transformation matrix being the identity matrix.

- Concerning the uniqueness of the Hermite normal form, we refer the reader to [20, Theorem II.3]; the Hermite normal form of a matrix $A$ is unique up to multiplication by units as long as $A$ is invertible.

We utilise the technique of the proof to obtain more examples of Hermite rings.

**Lemma 2.15** ([21, Theorem 3])**.** *Let $R$ be a commutative ring, then $R$ is a Hermite ring, iff*

$$\forall a, b \in S : \exists a_1, b_1, d \in R \text{ such that } a = a_1 d, b = b_1 d \text{ and } (a_1, b_1) = R. \qquad (2.1)$$

*Proof.*
"$\Leftarrow$" Let $R$ satisfy condition (2.1). Owing to the proof of Proposition 2.13, it suffices to demonstrate that every $1 \times 2$ matrix $\begin{pmatrix} a & b \end{pmatrix}$ permits a diagonal reduction. Therefore choose $a_1, b_1, d, s, t$ satisfying

$$a = a_1 d, b = b_1 d \text{ and } s a_1 + t b_1 = 1.$$

Put $Q = \begin{pmatrix} s & -b_1 \\ t & a_1 \end{pmatrix}$, then $Q$ is unimodular and $\begin{pmatrix} a & b \end{pmatrix} Q = \begin{pmatrix} d & 0 \end{pmatrix}$.

"$\Rightarrow$" To prove the other direction, let $R$ be a commutative Hermite ring. Let $a, b \in R$ be arbitrary. The definition of Hermite guarantees the existence of an invertible matrix $Q = \begin{pmatrix} s & -b_1 \\ t & a_1 \end{pmatrix}$, such that

$$\begin{pmatrix} a & b \end{pmatrix} Q = \begin{pmatrix} d & 0 \end{pmatrix} \text{ for some } d \in R.$$

This yields $a b_1 = b a_1$, $sa + tb = d$. $Q$ being unimodular, we deduce $s a_1 + t b_1 = 1$. Now

$$s a_1 a + t b_1 a = a \text{ implies } s a_1 a + t a_1 b = a,$$

hence $a_1 d = a$. Analogously $b_1 d = b$, which shows that condition (2.1) is fulfilled.   □

As the notion will be useful in the sequel we briefly attend to the notion of a GCD-domain.

**Definition 2.16.** An integral domain $\mathcal{D}$ is called a *GCD-domain*, if

$$\forall a, b \in \mathcal{D} : \exists d \in \mathcal{D} : d|a \wedge d|b, \text{ such that } \forall d' \in \mathcal{D} : d'|a \wedge d'|b \Rightarrow d'|d.$$

The element $d$ being uniquely determined up to multiplication by units is then called the *greatest common divisor* of $a$ and $b$, denoted $\gcd(a, b)$.

**Remark 2.17.**

(i) It is not difficult to check that a Bézout integral domain $\mathcal{D}$ is a GCD-domain. In fact take arbitrary $a, b \in \mathcal{D}$. There exists $d \in \mathcal{D}$ such that $(a) + (b) = (d)$, which implies the existence of $r, s \in \mathcal{D}$ such that $ar + bs = d$. On the other hand $(a), (b) \subseteq (d)$ and hence $(d)|(a), (d)|(b)$, which leads to $d|a$ and $d|b$. If for any $d' \in \mathcal{D}$, we have $d'|a$ and $d'|b$, then there are $u, v \in \mathcal{D}$, such that $d'u = a$ and $d'v = b$. Substituting we arrive at $d'(ur + vs) = d$, showing that $d'|d$ and hence $d = \gcd(a, b)$.

(ii) Due to Theorem 2.9 and the latter remark, all principal ideal domains, commutative elementary divisor domains and Hermite integral domains are GCD-domains.

**Proposition 2.18** (cf. [17, Corollary 4.28]). *A commutative Bézout domain $\mathcal{D}$ is Hermite.*

*Proof.* We may obtain for all $a, b \in \mathcal{D}$ elements $a', b' \in \mathcal{D}$ such that

$$aa' + bb' = \gcd(a, b) =: d.$$

There exist $a_1, b_1$ satisfying $a_1 d = a$ and $b_1 d = b$, which leads to $a'a_1 + b'b_1 = 1$. Invoking the previous Lemma 2.15 evidences the claim. $\qquad\square$

**Remark 2.19.** The question, whether a commutative Bézout domain, i.e. a Hermite domain, already constitutes an elementary divisor domain has not been answered yet, though several conditions are known.[22][23][19]

After discussing the Hermite normal form, we finally proceed to the Smith Normal Form, which we need to settle Theorem 2.9.

We need the following auxiliary

**Lemma 2.20** ([20, Corollary II.1] or [24, §21]). *Let $\mathcal{D}$ denote a Bézout integral domain. Choose arbitrary $b_1, \ldots, b_n \in \mathcal{D}$. There exists a matrix $Q$ with first row $(b_1, \ldots, b_n)$ and $\det(Q) = \gcd(b_1, \ldots, b_n)$.*

*Proof.* The claim is trivial for $n = 1$. Considering the case $n = 2$, for $b_1, b_2 \in \mathcal{D}$ we find $r, t \in \mathcal{D}$ such that $\gcd(b_1, b_2) = b_1 r + b_2 t$. Thus the matrix $\begin{pmatrix} b_1 & b_2 \\ -t & r \end{pmatrix}$ has determinant $\gcd(b_1, b_2)$.

Suppose the statement holds for some $n - 1$. Set $d_{n-1} := \gcd(b_1, \ldots, b_{n-1})$, and $d_n = \gcd(b_1, \ldots, b_n) = \gcd(d_{n-1}, b_n)$. We find $\rho, \sigma \in \mathcal{D}$ with $\rho d_{n-1} - \sigma b_n = d_n$. Let $D_{n-1}$ be a matrix obtained by virtue of the induction hypothesis and set

$$
D_n = \left( \begin{array}{cccc|c}
 & & & & b_n \\
 & D_{n-1} & & & 0 \\
 & & & & \vdots \\
 & & & & 0 \\
\hline
\frac{b_1 \sigma}{d_{n-1}} & \frac{b_2 \sigma}{d_{n-1}} & \cdots & \frac{b_{n-1} \sigma}{d_{n-1}} & \rho
\end{array} \right).
$$

Now $D_n$ has $(b_1, \ldots, b_n)$ as top row. Using Laplace expansion on the last column, we have

$$
\det(D_n) = \rho \det(D_{n-1}) + (-1)^{n+1} b_n \det(E_{n-1}),
$$

where $E_{n-1}$ is the submatrix, that arises from omitting the first row and last column in $D_n$. One has

$$
d_{n-1} E_{n-1} = \begin{pmatrix}
0 & d_{n-1} & 0 & \cdots & 0 \\
0 & & d_{n-1} & & \vdots \\
\vdots & & & \ddots & 0 \\
0 & & & & d_{n-1} \\
\sigma & 0 & 0 & \cdots & 0
\end{pmatrix} D_{n-1}.
$$

As $\det(D_{n-1}) = d_{n-1}$ by the induction hypothesis, we compute

$$
d_{n-1}^{n-1} \det(E_{n-1}) = (-1)^{(n-1)+1} \sigma d_{n-1}^{n-2} \det(D_{n-1}) = (-1)^n \sigma d_{n-1}^{n-1},
$$

accordingly $\det(E_{n-1}) = (-1)^n \sigma$ and $\det(D_n) = \rho d_{n-1} - \sigma b_n = d_n \in \mathcal{D}^*$. This evidences that $D_n$ is unimodular. $\qquad \square$

Observe that it is just a matter of notation to obtain the result for columns instead of rows.

**Definition 2.21.** Let $R$ be a commutative ring, $A$ an $m \times n$ matrix over $R$. A *submatrix* of $A$ is any square matrix of $A$, that results from deleting columns or rows of $A$. We define the *rank* of $A$ to be the size of the largest submatrix with non-vanishing determinant.

Finally we have constructed the theoretical framework to prove

**Theorem 2.22** ([24, Theorem 26.2]). *A matrix over a principal ideal domain $\mathcal{D}$ admits a Smith Normal Form, i.e. a diagonal reduction.*[1]

*Proof.* Let $(a_{ij})$ be an $m \times n$ matrix of rank $r$. We may assume that $(a_{ij})$ possesses non-zero elements. As $\mathcal{D}$ is in particular a unique factorisation domain, the function $\Omega$ mapping an element to its total number (counting multiplicity) of prime factors is well-defined. We provide a step-by-step instruction on how to reduce a matrix to its Smith Normal Form:

(i) Use elementary operations on $(a_{ij})$ to shift an $r \times r$ submatrix with non-zero determinant to the left upper corner.

(ii) If $0 \neq a_{11}|a_{1j}$ for $1 \leq j \leq n$, set $(b_{ij}) = (a_{ij})$ go to step (iv).

(iii) Since $\mathcal{D}$ is in particular a commutative Bézout domain, we find $k_j \in \mathcal{D}, j \in \{1, \ldots, n\}$, such that $\sum_{j=1}^{n} k_j a_{1j} = \gcd(a_{11}, \ldots, a_{1n}) =: d$. Invoking the former Lemma 2.20 yields an $n \times n$ matrix $Q$ with first column $(k_1, \ldots, k_n)$ and determinant $\gcd(k_1, \ldots, k_n) = 1$.[2] Then $(b_{ij}) := (a_{ij})Q$ satisfies $0 \neq d = b_{11}|b_{1j}$ for $1 \leq j \leq n$, as the $b_{1j}$'s are linear combinations of the $a_{1j}$'s. Furthermore we have $\Omega(b_{11}) < \Omega(a_{11})$.[3]

(iv) If $0 \neq b_{11}|b_{i1}$ for $1 \leq i \leq m$, skip this step, else:
Repeat the previous step with columns and rows interchanged.

(v) Iterate steps (ii) to (iv) until $b_{11}$ divides all other elements within its column and row. The procedure terminates as the number of prime factors of $b_{11}$ are reduced in every step.

(vi) Apply elementary operations to produce zeros in the column below and the row beneath $b_{11}$.

(vii) Use the above steps on the remaining $(m-1) \times (n-1)$ matrix, leaving the first row and column unchanged. Inductively we arrive at a matrix $A \sim (a_{ij})$ of the form $A = \text{diag}(d_1, \ldots, d_r, M)$ for a matrix $M$ of appropriate size. We find $M = \mathbf{0}$, since otherwise a non-zero element in $M$ could be permuted to position $(r+1, r+1)$, which would yield a submatrix of size $r+1$ with non-vanishing determinant - contrary to the assumption imposed on the rank.

---

[1]The elementary divisors are uniquely determined up to associated elements.
[2]Let $\sum_{j=1}^{n} k_j a_{1j} = d$, then there exist $r_j \in \mathcal{D}$, such that $d r_j = a_j$. Hence $\sum_{j=1}^{n} k_j d r_j = d$ entailing $\sum_{j=1}^{n} k_j r_j = 1$ and hence $\gcd(k_1, \ldots, k_n) = 1$.
[3]Evidently $\Omega(b_{11}) > \Omega(a_{11})$ is impossible in view of $b_{11} = \gcd(a_{11}, \ldots, a_{1n})$. Suppose thus $\Omega(b_{11}) = \Omega(a_{11})$, then $0 \neq b_{11} = a_{11}$ and $a_{11}|a_{1j}$ for $1 \leq j \leq n$ - this would have caused the algorithm to skip from step (ii) to step (iv).

(viii) By adding columns one derives

$$\mathrm{diag}(d_1,\ldots,d_r) \sim \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ d_2 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ d_r & 0 & \cdots & d_r \end{pmatrix}.$$

Owing to Lemma 2.20 we find $(c_{ij}) \sim \mathrm{diag}(d_1,\ldots,d_r)$ with $c_{11} = \gcd(d_1,\ldots,d_r)$.

(ix) Repeat the last step for the remaining $(r-1) \times (r-1)$-matrix, leaving the first row and column unchanged, until $d_i|d_{i+1}$ for $1 \le i \le r-1$.

$\square$

**Remark 2.23.**

(i) Generalisations include (almost) simultaneous transformation of a set of matrices to their Smith Normal Form [25] and relaxation of principal ideal domain to special *Prüfer domains* [26]; A ring constitutes a Prüfer domain if there are no zero-divisors and every ideal is invertible, thus essentially being a Dedekind domain with the Noetherian condition dropped - or as discussed in Chapter 4 a semihereditary integral domain.

(ii) For a constructive algorithm computing the Smith Normal Form, where $\mathcal{D}$ is a ring of algebraic integers, see Cohen [27, Algorithm 4.4]. For a discussion of the non-commutative case see [28].

(iii) The theorem shows that two $m \times n$ matrices over a commutative principal ideal domain are equivalent, if and only if they can be reduced to the same Smith Normal Form. If we replace the notion of ring by field, the Smith Normal Form may be multiplied by elementary matrices to derive the form

$$\mathrm{diag}(\underbrace{1,\ldots,1}_{r\text{-times}}, \mathbf{0}_{(m-r)\times(n-r)}).$$

The number $r$ of non-vanishing entries, which corresponds to the rank of the matrix $A$, is invariant for every equivalence class $[A]_\sim$.

## 2.4 Application of normal forms in module theory

Though not our primary topic we are too close to proving the invariant factor decomposition for finitely generated modules over principal ideal domains to refrain from doing so. We chose to skip the canonical procedure of introducing presentations, relations and generators of finitely generated modules, but instead abbreviate by using

**Lemma 2.24** ([29, §20 Lemma 1.1]). *Let $\mathcal{F}$ be a free module of finite rank $n$ over some principal ideal domain $\mathcal{D}$. Then every submodule $\mathcal{U} \subseteq \mathcal{F}$ is free, its rank not exceeding $n$.*

*Proof.* Without loss of generality $\mathcal{F} \cong \mathcal{D}^n$. We may use induction on $n \geq 0$, the case $n = 0$ holds true trivially. Identifying $\mathcal{U}$ with its image in $\mathcal{D}^n$ consider the module epimorphism $\pi$ being the canonical projection of $\mathcal{U}$ onto its first component. If $\pi \equiv 0$, then $\mathcal{U} = \ker(\pi) \subseteq \mathcal{D}^{n-1}$ and we are done invoking induction. If $\pi \not\equiv 0$, $\pi(\mathcal{U})$ is an non-zero ideal in $\mathcal{D}$, using that $\mathcal{D}$ is a principal ideal domain we have $\pi(\mathcal{U}) = a\mathcal{D}$ for a suitable $a \in \mathcal{D} \setminus \{0\}$. We find $k \in \mathcal{U}$ with $\pi(k) = a$. For any $x \in \mathcal{U}$ we may write $\pi(x) = ba$, $b \in \mathcal{D}$. Now $x = bk + (x - bk)$, the last summand being contained in $\ker(\pi)$, we see $\mathcal{U} = k\mathcal{D} + \ker(\pi)$. A brief calculation evidences $k\mathcal{D} \cap \ker(\pi) = \{0\}$, implying that $\mathcal{U} = k\mathcal{D} \oplus \ker(\pi)$. By induction hypothesis the claim holds for $n - 1$. As $\ker(\pi) \subseteq \{(0, x_2, \ldots, x_n) | x_i \in \mathcal{D}\}$ injects into $\mathcal{D}^{n-1}$, $\ker(\pi)$ is free of rank at most $n - 1$, say $m$. We obtain

$$\mathcal{U} = \ker(\pi) \oplus k\mathcal{D} \cong \mathcal{D}^m \oplus \mathcal{D} \cong \mathcal{D}^{m+1}.$$

$\square$

We are now able to state and prove the invariant factor decomposition, which we will generalise to Dedekind domains in Chapter 4.

**Theorem 2.25** ([29, §21 Theorem 1.1]). *Let $M$ be a finitely generated module over a principal ideal domain $\mathcal{D}$, then*

$$M \cong \mathcal{D}^s \oplus \bigoplus_{i=1}^{r} \mathcal{D}/d_i\mathcal{D}$$

*with non-zero $d_i$ in $\mathcal{D} \setminus \mathcal{D}^*$ fulfilling $d_i | d_{i+1}$ for $1 \leq i \leq r - 1$.*

*Proof.* It is well-known that $M \cong \mathcal{D}^n/\mathcal{U}$ for a submodule $\mathcal{U}$ of the free module $\mathcal{D}^n$ with suitable $n \in \mathbb{N}$. Using the previous Lemma 2.24, we find a basis $\mathbf{a}_1, \ldots, \mathbf{a}_m$ of $\mathcal{U}$ with $m \leq n$. Denote by $A$ the $m \times n$ matrix obtained by taking the $\mathbf{a_i}$'s as rows. The Smith Normal Form guarantees the existence of $P \in \mathrm{GL}_m(\mathcal{D}), Q \in \mathrm{GL}_n(\mathcal{D})$ and $k \in \mathbb{N}$, such that

$$PAQ = \mathrm{diag}(d_1, \ldots, d_k, \mathbf{0}_{(m-k)\times(n-k)}), \text{ where } d_i \mid d_{i+1} \text{ for } 1 \leq i \leq k - 1.$$

Since $\mathcal{U} = A\mathcal{D}^n = A.\mathbf{e}_1 \oplus \cdots \oplus A.\mathbf{e}_n$ we calculate

$$M \cong \mathcal{D}^n/\mathcal{U} = \mathcal{D}^n/A\mathcal{D}^n \cong P\mathcal{D}^n/PAQ\mathcal{D}^n,$$

the last isomorphy being valid as a multiplication by invertible matrices $P, Q$ merely constitutes a change of bases. Let $r$ be the number of those $d_i$'s, which are not units.

Then

$$M \cong \mathcal{D}^n \big/ \mathrm{diag}(d_1, \ldots, d_k, \mathbf{0}_{(m-k)\times(n-k)})\mathcal{D}^n \cong \mathcal{D}^{n-k} \oplus \bigoplus_{i=1}^{r} \mathcal{D} \big/ d_i \mathcal{D},$$

where the $d_i$'s, that are units, do not contribute to the last term, as in this case $\mathcal{D}^n \big/ d_i \mathcal{D}^n \cong 0$. The rank of the free part $\mathcal{D}^{n-k}$ stems from the number of zero columns.

$\square$

**Remark 2.26.** By setting $\mathcal{D} = \mathbb{Z}$ the *fundamental theorem for finitely generated abelian groups* follows as an immediate corollary.

## 2.5 Inheritance

A facile yet rewarding approach to the unit sum number problem for matrices is inheritance with respect to the matrix size; concretely let $R$ be a unital, commutative ring and suppose $u\big(\mathrm{Mat}_n(R)\big) = 2$, can the same be deduced for matrix rings of higher or lower size?

The next results are of utter usefulness to the *complete* determination of $u\big(\mathrm{Mat}_n(R)\big)$ for all sizes $n > 1$.

**Proposition 2.27** ([2, Theorem 12]). *Assume that the set $R \setminus R^*$ is 2-good, then $u\big(\mathrm{Mat}_n(R)\big) = 2$, $\forall n > 1$.*

*Proof.* Fix $n > 1$ and take $A \in \mathrm{Mat}_n(R)$. Suppose every element on the diagonal of $A$ is 2-good, then $A$ admits a decomposition into an upper triangular and a lower triangular matrix respectively - both invertible as the diagonal entries are units. Therefore assume that the number of units on the diagonal is $\geq 1$, say $r$. By elementary row and column operations we swap a unit to position $(1,1)$, multiplying appropriately we can clear out the column below and the row beneath the unit, so that

$$A \sim \begin{pmatrix} 1 & 0 \ldots 0 \\ 0 & \\ \vdots & B \\ 0 & \end{pmatrix}.$$

Iterating the process $A$ becomes equivalent to a matrix of the form $A' = \begin{pmatrix} I_r & 0 \\ 0 & B \end{pmatrix}$, where $B \in \mathrm{Mat}_{n-r}(R)$ has only 2-good elements for its diagonal entries. If $r = n$ or $n = 2$ then Proposition 2.5 yields that $A$ is 2-good. Let $n > 2$ and $r = 1$, a

decomposition of $A'$ is given by

$$U = \begin{pmatrix} 1 & 1 & 0 & \ldots & 0 \\ 1 & 0 & b_{1,2} & \ldots & b_{1,n} \\ 0 & 0 & u_2 & \ldots & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdot & \ldots & u_{n-1} \end{pmatrix} \text{ and } V = \begin{pmatrix} 0 & -1 & 0 & \ldots & 0 \\ -1 & b_{1,1} & 0 & \ldots & 0 \\ 0 & b_{2,1} & v_2 & \ldots & \cdot \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & b_{n-1,1} & \cdot & \ldots & v_{n-1} \end{pmatrix},$$

where we decomposed the diagonal elements $b_{ii} = u_i + v_i$, with $u_i, v_i \in R^*$. Laplace expansion along the first columns of $U$ and $V$ evidences $U, V \in \mathrm{GL}_n(R)$. Whence we are left with $r \geq 2$. Now $I_r = P_1 + Q_1$ due to Proposition 2.5 and $B = P_2 + Q_2$ as the diagonal entries of $B$ are 2-good with $P_i, Q_i$ invertible, as discussed at the beginning of the proof. Finally

$$\begin{pmatrix} P_1 & 0 \\ 0 & P_2 \end{pmatrix}, \begin{pmatrix} Q_1 & 0 \\ 0 & Q_2 \end{pmatrix} \in \mathrm{GL}_n(R),$$

due to Remark 2.2(ii), which demonstrates that $A$ is 2-good. $\square$

**Corollary 2.28** ([2, Corollary 14]). *If $u\big(\mathrm{Mat}_n(R)\big) = 2$ for some $n \in \mathbb{N}$, then $u\big(\mathrm{Mat}_{nk}(R)\big) = 2$, $\forall k \in \mathbb{N}$.*

*Proof.* Apply the prior Proposition 2.27 to $\mathrm{Mat}_k\big(\mathrm{Mat}_n(R)\big)$. $\square$

**Proposition 2.29** ([12, Proposition 8]). *If $m, n > 1$, such that*

$$u\big(\mathrm{Mat}_n(R)\big) = u\big(\mathrm{Mat}_m(R)\big) = 2,$$

*then $u\big(\mathrm{Mat}_{n+m}(R)\big) = 2$*

*Proof.* Take a matrix $M \in \mathrm{Mat}_{n+m}(R)$ and use a block representation $M = \begin{pmatrix} A & B \\ C & D \end{pmatrix}$, where $A \in \mathrm{Mat}_n(R), D \in \mathrm{Mat}_m(R)$, the blocks $B$ and $C$ of appropriate size. By assumption there are $A_1, A_2 \in \mathrm{GL}_n(R)$ and $D_1, D_2 \in \mathrm{GL}_n(R)$ such that $A = A_1 + A_2$ and $D = D_1 + D_2$. Now

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} A_1 & B \\ 0 & D_1 \end{pmatrix} + \begin{pmatrix} A_2 & 0 \\ C & D_2 \end{pmatrix}$$

and Laplace expansion shows that the summands are invertible. $\square$

**Remark 2.30.** We conclude the chapter by reviewing the main results regarding the unit sum number of matrix rings over certain rings.

- In Corollary 2.6 we proved that no matter the ring $S$, the matrix ring $\mathrm{Mat}_n(S)$ of size $n > 1$ is 3-good. Therefore the unit sum number - i.e. the minimal $k \in \mathbb{N}_{>1}$ by which every matrix over $S$ may be decomposed into a sum of exactly $k$ invertible matrices - of a matrix ring is either 2 or 3.

- A matrix ring over a principal ideal domain has unit sum number 2, as all principal ideal domains constitute an elementary divisor domain due to Proposition 2.9.

- Although not apparent from the literature, in view of $2\mathbb{Z} + 3\mathbb{Z} = \mathbb{Z}$ it is evident that the results just proved in Corollary 2.28 and Proposition 2.29 lead to:

$$u\big(\operatorname{Mat}_k(R)\big) = 2 \ \forall k > 1 \Leftrightarrow u\big(\operatorname{Mat}_2(R)\big) = u\big(\operatorname{Mat}_3(R)\big) = 2,$$

for any commutative ring $R$.

# Semilocal rings

A not necessarily commutative ring $S$ is called semilocal, if $S/\mathrm{Jac}(S)$ is semisimple.[1] Recall that an $S$-module $M$ is semisimple, if it can be written as (direct) sum of simple modules, i.e. modules featuring solely $0$ and $M$ as submodules. The ring $S$ is then called semisimple, if it is semisimple as an $S$-module. Note that every Artinian ring and a fortiori every local ring is semilocal.[2] We are able to discuss the unit sum number problem for semilocal rings exhaustively.

## 3.0.1 Artin-Wedderburn structure theorem

One of the ingredients needed to settle the unit sum number problem, is the Artin-Wedderburn structure theorem on semisimple rings, which we will develop in the sequel. The results and proofs in this section are taken from a lecture course by Joachim Mahnkopf about representation theory [31].

**Proposition 3.1.** *Let $S$ be a semisimple ring. There are only finitely many isomorphism classes of simple $S$-modules. If $\{M_\alpha\}_{\alpha \in A}$ denotes a system of representatives of isomorphism classes of simple $S$-modules, we have*

$$S \cong \bigoplus_{\alpha \in A} n_\alpha M_\alpha := \bigoplus_{\alpha \in A} \bigoplus_{i=1}^{n_\alpha} M_\alpha.$$

*Proof.* Let $M$ be a simple $S$-module. First we deduce the existence of minimal left ideals $I_1, \ldots, I_m$ of $S$, such that $S = \bigoplus_{i=1}^{m} I_i$.
As $S$ is semisimple, we have $S \cong \bigoplus_{i \in I} I_i$, where the $I_i \leq S$ are simple submodules of $S$. We find $s_i \in I_i$, almost all of them equal to zero, such that $1 = \sum_{i \in I} s_i$. Suppose

---

[1]We do not distinguish between left- and right-semisimple, as one property implies the other.
[2]cf. [30, (20.3)]

the index set $I$ is infinite, then there is a $i_0 \in I$, such that $s_{i_0} = 0$. Choose $0 \neq s \in I_{i_0}$, then

$$s = s \cdot 1 = \sum_{i \in I} s s_i = \sum_{i \in I \setminus \{i_0\}} s s_i.$$

Thus $I_{i_0} \cap \bigoplus_{i \in I \setminus \{i_0\}} I_i \neq \{0\}$, contradicting the directness of $\bigoplus_{i \in I} I_i$ and hence $I$ is finite.

Next we show that, $M$ is isomorphic to one of the $I_i$. From this we conclude that $\{I_1, \ldots, I_m\}$ contains a system of representatives of isomorphism classes of simple $S$-modules, finishing the proof.
As $M$ is simple, we find $m \in M$, such that $< m >= M$. Defining as usual $\mathrm{Ann}(m) = \{s \in S : sm = 0\}$, we find a non-trivial homomorphism

$$\pi : S \to \mathop{S}\!/\!\mathrm{Ann}(m) \xrightarrow{\cong} M.$$

Hence

$$\{0\} \neq \mathrm{Hom}_S(S, M) = \mathrm{Hom}_S\Big(\bigoplus_{i=1}^m I_i, M\Big) = \bigoplus_{i=1}^m \mathrm{Hom}_S(I_i, M).$$

We find $i_0 \in \{1, \ldots, m\}$ and a non-trivial homomorphism $\varphi : I_{i_0} \to M$. As $I_{i_0}$ and $M$ are simple, $\varphi$ constitutes an isomorphism. As each simple $S$-module is isomorphic to one of the ideals $I_1, \ldots, I_m$, the claim follows.  $\square$

Next we provide a part of the well-known Schur-Lemma.

**Lemma 3.2.** *The endomorphism ring* $\mathrm{End}_S(M)$ *of a simple $S$-module $M$ is a division ring.*

*Proof.* Let $0 \neq f \in \mathrm{End}_S(M)$. As $\ker(f) \leq M$ and $\mathrm{im}(f) \leq M$, we can only have $\ker(f) = 0$ as $0 \neq f$, showing that $f$ is injective. Analogously $\mathrm{im}(f) = M$, showing that $f$ is indeed invertible, whence $\mathrm{End}_S(M)$ constitutes a division ring.  $\square$

We collect some rules for calculations with endomorphism rings.

**Lemma 3.3.** *Let $S$ be an arbitrary ring, then*

(i) $\mathrm{End}_S(S) \cong S^{\mathrm{op}}$

(ii) $\mathrm{Mat}_n(S)^{\mathrm{op}} \cong \mathrm{Mat}_n(S^{\mathrm{op}})$

*Proof.*

(i) Consider the following map

$$\theta : S^{\mathrm{op}} \to \mathrm{End}_S(S) \text{ given by } s \mapsto \{\theta_s : m \mapsto ms\}.$$

Clearly $\theta_s \in \mathrm{End}_S(S)$, some unproblematic calculations evidence that $\theta$ is a bijective homomorphism of rings.

(ii) Define
$$\mu : \mathrm{Mat}_n(S)^{\mathrm{op}} \to \mathrm{Mat}_n(S^{\mathrm{op}}) \text{ by } A \mapsto A^{\mathsf{T}}.$$

It is readily observed that $\mu$ fulfils the requirements.

$\square$

**Proposition 3.4** ([32, §8 Lemma 2.7])**.** *Let $M$ be an $S$-module and $E = \mathrm{End}_S(M)$, then*
$$\mathrm{End}_S(nM) \cong \mathrm{Mat}_n(E) \quad \forall n \in \mathbb{N}.$$

*Proof.* Define $\Phi : \mathrm{Mat}_n(E) \to \mathrm{End}_S(nM)$ via $(\varphi_{ij}) \mapsto \varphi$, where $\varphi : nM \to nM$ is given by
$$\varphi(x_1, \dots, x_n) = \Big( \sum_{j=1}^{n} \varphi_{1j}(x_j), \dots, \sum_{j=1}^{n} \varphi_{nj}(x_j) \Big).$$

By using linearity of the sums within the components, one checks that $\varphi$ is $S$-linear and $\Phi$ is a homomorphism of rings.

To prove the opposite direction define $\Psi : \mathrm{End}_S(nM) \to \mathrm{Mat}_n(E)$ via $\psi \mapsto (\psi_{ij})$, where the $\psi_{ij} : M \to M$ are given by
$$\big( \psi_{1j}(x), \dots, \psi_{nj}(x) \big) = \psi(0, \dots, 0, \underbrace{x}_{\text{j-th pos.}}, 0, \dots, 0).$$

Again one checks $\psi_{ij} \in \mathrm{End}_S(M)$ with $1 \leq i, j \leq n$ and $\Psi$ is a ring homomorphism. Finally keeping an account of indices and summation one verifies $\Psi = \Phi^{-1}$. $\square$

Clearly, the proposition is a generalisation of the well-known result from linear algebra, which establishes the isomorphy of linear maps between vector spaces and matrices over the base field.

Based on Wedderburn's work [33] in 1907, Artin in 1927 developed the famous structure theorem for semisimple rings, which serves as a cornerstone in the theory of non-commutative rings. A comprehensive synopsis may be found in [30].

**Theorem 3.5** (Artin-Wedderburn)**.** *Let $S$ be a semisimple ring. Then there exist $n_1, \dots, n_r \in \mathbb{N}$ and division rings $D_1, \dots, D_r$ such that*
$$S \cong \bigoplus_{i=1}^{r} \mathrm{Mat}_{n_i}(D_i).$$

*Proof.* Since $S$ is semisimple we find $n_1, \ldots, n_r \in \mathbb{N}$ and pairwise non-isomorphic, simple $S$-modules $M_i \leq S$ such that $S = \bigoplus_{i=1}^{r} n_i M_i$. Setting $D_i = \mathrm{End}_S(M_i)$, Lemma 3.2 indicates that $D_i$ is a division ring, as $M_i$ is simple. Employing the former lemmata we compute

$$
\begin{aligned}
S^{\mathrm{op}} &\cong \mathrm{End}_S(S) \\
&= \mathrm{End}_S \Big( \bigoplus_{i=1}^{r} n_i M_i \Big) = \mathrm{Hom}_S \Big( \bigoplus_{i=1}^{r} n_i M_i, \bigoplus_{i=1}^{r} n_i M_i \Big) \\
&= \bigoplus_{1 \leq i,j \leq r} \mathrm{Hom}_S(n_i M_i, n_j M_j) = \bigoplus_{i=1}^{r} \mathrm{Hom}_S(n_i M_i, n_i M_i) \\
&= \bigoplus_{i=1}^{r} \mathrm{End}_S(n_i M_i) \cong \bigoplus_{i=1}^{r} \mathrm{Mat}_{n_i}(D_i).
\end{aligned}
$$

The equality in the third line of the equation holds since for $i \neq j$ the $M_i$ and $M_j$ are non-isomorphic due to assumption and being simple they warrant $\mathrm{Hom}_S(n_i M_i, n_j M_j) = \{0\}$. Finally

$$
S = \bigoplus_{i=1}^{r} \mathrm{Mat}_{n_i}(D_i)^{\mathrm{op}} = \bigoplus_{i=1}^{r} \mathrm{Mat}_{n_i}(D_i^{\mathrm{op}}).
$$

As $D_i$ is a division ring, so is $D_i^{\mathrm{op}}$; this completes the proof.     $\square$

### 3.0.2  Main result and Zelinsky's theorem

The main result of this section is

**Theorem 3.6** ([12, Lemma 2]). *For a semilocal ring $S$*

(i) $u(S) = 2$, *if there is no factor ring isomorphic to $\mathbb{F}_2$.*

(ii) $u(S) = \omega$, *if there exists exactly one factor isomorphic to $\mathbb{F}_2$.*

(iii) $u(S) = \infty$, *if $\mathbb{F}_2 \oplus \mathbb{F}_2$ is a factor.*

In order to prove the first claim, we use one of the earliest results concerning sums of units by Zelinsky about matrix rings together with the Artin-Wedderburn structure theorem on semisimple rings.

For the purpose of proving Theorem 3.6 it suffices to state a simple, finite dimensional version of Zelinsky's theorem.

**Theorem 3.7** ([1]). *Let $V$ be a vector space of dimension $n$ over a division ring $D$. Then*

$$u\big(\operatorname{Mat}_n(D)\big) = \begin{cases} \omega & \text{if } D = \mathbb{F}_2 \text{ and } n = 1 \\ 2 & \text{else.} \end{cases}$$

*Proof.* Employing Lemmata 3.3 and 3.4 we have

$$\operatorname{Mat}_n(D) \cong \operatorname{Mat}_n(D)^{\operatorname{op}} \cong \operatorname{Mat}_n(D^{\operatorname{op}})$$
$$\cong \operatorname{Mat}_n\big(\operatorname{End}_D(D)\big) \cong \operatorname{End}_D(nD) \cong \operatorname{End}_D(V).$$

Lemma 1.3 (i) assures that $\operatorname{End}_D(V)$ is $k$-good, if and only if $\operatorname{Mat}_n(D)$ is $k$-good for some $k \in \mathbb{N}_{>1} \cup \{\omega\}$. Hence it suffices to prove the statement for the endomorphism ring of $V$ over $D$.

Assume first that $D = \mathbb{F}_2$ and $n = 1$, then also $V \cong \mathbb{F}_2$. Lemma 3.3 yields $\operatorname{End}_{\mathbb{F}_2}(\mathbb{F}_2) \cong \mathbb{F}_2$. Therefore $u\big(\operatorname{End}_D(V)\big) = u(\mathbb{F}_2) = \omega$ due to Lemma 1.1, which evidences the first part of the claim.

For the general case, where $D \neq \mathbb{F}_2$ or $n \neq 1$, take an arbitrary $\alpha \in \operatorname{End}_D(V)$. Choose complements $M, N$ such that $V = \ker\alpha \oplus M$ and $V = \operatorname{im}\alpha \oplus N$. In particular there exists an isomorphism $\varphi : \ker\alpha \to N$. We start by decomposing the identity map $\operatorname{id} := \operatorname{id}_{|\operatorname{im}\alpha}$ into a sum of two isomorphisms $\sigma$ and $\operatorname{id} - \sigma$.

- If the space is one-dimensional $\sigma$ may be represented as matrix $(s)$, where $s$ is a single element of $D \setminus \{0, 1\}$. Clearly the exceptional case $D = \mathbb{F}_2$ prevents this choice of $(s)$. (cf. Lemma 1.1)

- If $\dim(\operatorname{im}\alpha) = 2$, set $\operatorname{id} = \sigma + (\operatorname{id} - \sigma)$, where $\sigma = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$.

- If $\dim(\operatorname{im}\alpha) = 3$, set $\operatorname{id} = \sigma + (\operatorname{id} - \sigma)$, where $\sigma = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$. By multiplication
  with elementary matrices it is easy to show that the matrix corresponding to $(\operatorname{id} - \sigma)$ can be transformed to the identity matrix and is therefore invertible.

- If $\dim(\operatorname{im}\alpha) > 3$ decompose $\operatorname{im}\alpha$ into two- and three-dimensional subspaces and let $\sigma$ be the direct product of the corresponding two- and three-dimensional automorphisms.

We acquire $\sigma \in \operatorname{Aut}_D(\operatorname{im}\alpha)$ such that $(\operatorname{id} - \sigma) \in \operatorname{Aut}_D(\operatorname{im}\alpha)$. Consider the following mappings

$$\beta = \varphi \oplus (\sigma \circ \alpha_{|M}) \text{ and } \gamma = -\varphi \oplus \big((\operatorname{id} - \sigma) \circ \alpha_{|M}\big)$$

in $\operatorname{Aut}_D(V)$. We obtain $\beta + \gamma = 0 \oplus \alpha_{|M'} = \alpha$ and have thus shown that an arbitrary endomorphism $\alpha$ is 2-good, whence the same applies to $\operatorname{End}_D(V)$. $\qquad\square$

Eventually we have gathered enough results to prove the main theorem of the section.

*Proof of Theorem 3.6.* Let $S$ be a semilocal ring. Due to Lemma 1.3(ii) proving the assertions for $S/\mathrm{Jac}(S)$ suffices. Now $S/\mathrm{Jac}(S)$ is semisimple. Invoking Theorem 3.5 and adopting its notation, we obtain the representation

$$S/\mathrm{Jac}(S) \cong \bigoplus_{i=1}^{r} \mathrm{Mat}_{n_i}(D_i).$$

It is easy to see that

$$\mathrm{Mat}_{n_i}(D_i) \cong \mathbb{F}_2 \text{ if and only if } n_i = 1 \text{ and } D_i \cong \mathbb{F}_2.^3$$

(i) As $\mathbb{F}_2$ is not contained in the representation, Lemma 1.3(iii) together with Theorem 3.7 yields $u\big(S/\mathrm{Jac}(S)\big) = 2$.

(ii) Suppose the representation contains exactly one factor isomorphic to $\mathbb{F}_2$, i.e.

$$S/\mathrm{Jac}(S) \cong \mathbb{F}_2 \oplus \bigoplus_{i=1}^{r-1} \mathrm{Mat}_{n_i}(D_i).$$

As every $\mathrm{Mat}_{n_i}(D_i)$ is 2-good, so is $\bigoplus_{i=1}^{r-1} \mathrm{Mat}_{n_i}(D_i)$ due to Lemma 1.3(iii). Consider the projection

$$\varphi : \mathbb{F}_2 \oplus \bigoplus_{i=1}^{r-1} \mathrm{Mat}_{n_i}(D_i) \twoheadrightarrow \mathbb{F}_2.$$

If $S/\mathrm{Jac}(S)$ being the domain of the non-trivial epimorphism $\varphi$ would be $k$-good for some $k \in \mathbb{N}$, so would $\mathbb{F}_2$ - a contradiction to Lemma 1.3(i). On the other hand clearly every element in $S/\mathrm{Jac}(S)$ can be expressed as a sum of units, hence the ring is $\omega$-good.

(iii) The decomposition of $S/\mathrm{Jac}(S)$ takes the form

$$(\mathbb{F}_2 \oplus \mathbb{F}_2) \oplus \bigoplus_{i=1}^{r-2} \mathrm{Mat}_{n_i}(D_i) := V \oplus W.$$

Take an arbitrary $x \in W$, then $\big((1,0), x\big) \in V \oplus W$ cannot be expressed as sum of units in $V \oplus W$. Thus $u\big(S/\mathrm{Jac}(S)\big) = \infty$.

<div style="text-align: right;">□</div>

---

[3] Direction "⇐" is trivial. For the other direction note that $n_i$ cannot be greater one, as in this case $\mathrm{Mat}_{n_i}(D)$ would contain zero-divisors. Clearly $\mathrm{Mat}_1(D) \cong \mathbb{F}_2$ implies $D \cong \mathbb{F}_2$.

# Connections to the class number

## 4.1 Formulation of the main result

This chapter is devoted to an astonishing connection between the unit sum number of matrix rings over Dedekind domains $\mathfrak{O}$ and the class number $h_{\mathfrak{O}}$ of $\mathfrak{O}$; Vámos and Wiegand established the following theorem providing more information on the subject of matrix rings featuring unit sum number 2.

**Theorem 4.1** ([15, Thm 4.7])**.** *Let $h$ denote the finite class number of a Dedekind domain $\mathfrak{O}$. Then $u\big(\mathrm{Mat}_n(\mathfrak{O})\big) = 2$ for all $n \geq 2h$.*

The strength of this result stems from the observation, that the unit sum number of the matrix rings may implicate a lower bound for the class number: suppose we would find some $n_0 \in \mathbb{N}$ with $u\big(\mathrm{Mat}_{n_0}(\mathfrak{O})\big) = 3$, then clearly $\frac{n_0}{2} < h$. Allowedly the instruments for proving, that a matrix is not representable as a sum of two invertible matrices are scarce; for a sufficient condition regarding some minor special cases see [15, Proposition 4.9]. A fully worked example of a matrix ring with unit sum number 3, will be given in Example 4.5 below.

For the rest of the chapter let $\mathfrak{O}$ denote a Dedekind domain with finite class number $h$. Levy [4] calls a matrix $A$ *decomposable*, if it is equivalent to the block diagonal sum of two matrices - otherwise *indecomposable*.

The strategy for proving Theorem 4.1 is to establish the subsequent three results.

**Lemma 4.2** ([15, Remark 1.2(2)])**.** *Every matrix is equivalent to a block diagonal sum of indecomposable matrices.*

*Proof.* Let $A$ be some $m \times n$ matrix, if $A$ is indecomposable there is nothing to prove. If $A$ is decomposable we find $A_1, A_2$ of size smaller than $A$, such that $A \sim \mathrm{diag}(A_1, A_2)$. Evidently this process of decomposing terminates as the sizes of the occurring blocks are reduced in every step. $\qquad\square$

Secondly we take advantage of a result by Levy introduced in 1972.

**Theorem 4.3** ([4, Theorem 2.2])**.** *Let $A$ be an indecomposable $m \times n$ matrix over $\mathfrak{O}$. Then both $m, n \leq h$.*

The theorem was generalised to certain Prüfer domains by Vámos and Wiegand [15] in 2011 - still we wish to follow the original exposure of Levy, which warrants more insight into the theory of Dedekind domains, especially the structure of finitely generated modules over Dedekind domains. Thirdly the afore mentioned authors found

**Proposition 4.4** ([15, Corollary 4.6])**.** *Let $B = \mathrm{diag}(B_1, \ldots, B_t)$ be a block diagonal sum of size $n \times n$, $n \geq 2$. Suppose the number of rows and columns of every $B_i$ is $\leq n/2$, then $B$ is the sum of two units.*

Putting these results together we arrive at the

*Proof of 4.1.* Let $\mathbb{N} \ni n \geq 2h$ and $A \in \mathrm{Mat}_n(\mathfrak{O})$ be arbitrary. Then $A$ is equivalent to a block diagonal sum, each block size not exceeding $h$ by Theorem 4.3 . Due to the matrix size the assumptions in the former proposition 4.4 are fulfilled, hence $A$ is the sum of two units, entailing $\mathrm{Mat}_n(\mathfrak{O}) = 2$. $\qquad\square$

As announced beforehand, we now turn to an example of a Dedekind domain $\mathcal{O}$, being a ring of algebraic integers, where the unit sum number $u\big(\mathrm{Mat}_2(\mathcal{O})\big)$ equals three.

**Example 4.5** ([15, Example 4.11] and cf. [12, Proposition 10])**.** Let $\mathcal{O} = \mathbb{Z}[\sqrt{-5}]$, it is well-known that the class number of $\mathcal{O}$ is 2. Now the previous theorem asserts that

$$u\big(\mathrm{Mat}_n(\mathcal{O})\big) = 2 \text{ for } n \geq 2h = 4.$$

Whether $u\big(\mathrm{Mat}_3(\mathcal{O})\big) = 2$ has not been determined yet, but we are able to derive $u\big(\mathrm{Mat}_2(\mathcal{O})\big) = 3$:
We start by showing that the equation

$$d = 3r + (2 + \sqrt{-5})t \text{ with } r, t \in \mathcal{O}$$

has no solutions for $d \in \{\pm 1, \pm 2\}$. Writing $r = r_1 + \sqrt{-5}\, r_2$ and $t = t_1 + \sqrt{-5}\, t_2$ with $r_i, t_i \in \mathbb{Z}$ we find

$$\begin{aligned} d &= 3(r_1 + \sqrt{-5}\, r_2) + (2 + \sqrt{-5})(t_1 + \sqrt{-5}\, t_2) \\ &= 3r_1 + 2t_1 - 5t_2 + \sqrt{-5}(3r_2 + t_1 + 2t_2), \end{aligned}$$

hence

$$3r_2 + t_1 + 2t_2 = 0 \text{ and } 3r_1 + 2t_1 - 5t_2 = d$$

and we obtain

$$-(6r_2 - 3r_1 + 9t_2) = d \in 3\mathbb{Z},$$

which yields $d \notin \{\pm 1, \pm 2\}$.

Set $\mathfrak{d} = 3\mathcal{O} + (2 + \sqrt{-5})\mathcal{O}$. The calculation just made asserts $\mathfrak{d} \neq \mathcal{O}$ as $\mathcal{O}^* = \{-1, 1\}$. Now suppose $\mathfrak{d}$ were principal, i.e. we find $g \in \mathcal{O}$ with $3\mathcal{O} + (2 + \sqrt{-5})\mathcal{O} = g\mathcal{O}$. Then $g \mid 3$ and $g \mid (2 + \sqrt{-5})$. Clearly $g$ is not associated to 3, as an element associated to 3 does not divide $2 + \sqrt{-5}$. And as 3 is prime in $\mathcal{O}$, we obtain $g \in \mathcal{O}^*$, contradicting $\mathfrak{d} \neq \mathcal{O}$.

Let

$$A = \begin{pmatrix} 3 & 0 \\ 2 + \sqrt{-5} & 0 \end{pmatrix}$$

and suppose there are $(u_{ij}) = U, V \in \mathrm{GL}_2(\mathcal{O})$ admitting $A = U + V$. Define

$$\begin{pmatrix} a_1 & 0 \\ a_2 & 0 \end{pmatrix} := U^{-1}A = I + U^{-1}V$$

and observe that $\mathfrak{d} = a_1\mathcal{O} + a_2\mathcal{O}$.[1] Now

$$U^{-1}A - I = U^{-1}V = \begin{pmatrix} a_1 - 1 & 0 \\ a_2 & -1 \end{pmatrix}$$

is invertible, therefore entailing $a_1 - 1 \in \mathcal{O}^* = \{\pm 1\}$, which leaves $a_1 \in \{0, 2\}$. As $a_1 = 2$ is ruled out by the calculation above, we must have $a_1 = 0$, implying that $\mathfrak{d}$ is a principal ideal - a contradiction showing that $A$ is not 2-good.

## 4.1.1 The unit sum number of block matrices

**Definition 4.6.**

(i) A *permutation matrix* is a square matrix having exactly one neutral element 1 in each row and in each column.

(ii) Let $P = (P_{ij})_{1 \leq i,j \leq n}$ be a permutation matrix, $X = (X_{ij})_{1 \leq i,j \leq n}$ a square matrix of the same size. Define the *meeting number* to be

$$m(P, X) := \#\{(i,j) \mid P_{ij} \neq 0 \wedge X_{ij} \neq 0\}.$$

We say $P$ avoids $X$, if $m(P, X) = 0$.

---

[1] We see $a_1 = 3u_{11} + (2 + \sqrt{-5})u_{12}$ and $a_2 = 3u_{21} + (2 + \sqrt{-5})u_{22}$, where $u_{11}u_{22} - u_{12}u_{21} = \det(U) \in \mathcal{O}^*$. It is evident that $a_1\mathcal{O} + a_2\mathcal{O} \subseteq \mathfrak{d}$. To see the other inclusion note that $(2 + \sqrt{-5})u_{22}, (2 + \sqrt{-5})u_{12} \in a_1\mathcal{O} + a_2\mathcal{O}$, thus also $(2 + \sqrt{-5})u_{11}u_{22}$ and $-(2 + \sqrt{-5})u_{12}u_{21}$ lie in $a_1\mathcal{O} + a_2\mathcal{O}$. This leads to $\det(U)(2 + \sqrt{-5}) \in a_1\mathcal{O} + a_2\mathcal{O}$. Therefore $(2 + \sqrt{-5})\mathcal{O} \subseteq a_1\mathcal{O} + a_2\mathcal{O}$. Analogously $3\mathcal{O} \subseteq a_1\mathcal{O} + a_2\mathcal{O}$ and hence $\mathfrak{d} \subseteq a_1\mathcal{O} + a_2\mathcal{O}$.

The upcoming two lemmata link the block size within a block diagonal sum to the presentation by two units.

**Lemma 4.7** ([15, Lemma 4.3]). *Let $S$ be a ring, $A \in \mathrm{Mat}_n(S)$ avoiding a permutation matrix $P \in \mathrm{Mat}_n(S)$. Then $A$ is the sum of two invertible matrices over $S$.*

*Proof.* We will prove the existence of $n \times n$ matrices $A_1, A_2$, which sum to $A$, such that $A_1 + P, A_2 - P \in \mathrm{GL}_n(S)$.

First assume $P$ is the identity matrix, which entails $\mathrm{diag}(A) = (0, \ldots, 0)$. Take $A_1$ to be the strict lower triangular part of $A$, the remaining entries filled with zeros. Utilizing the same construction on $A_2$ but employing the strict upper part of $A$, yields matrices $A_1, A_2$ satisfying the condition.

Now let $P$ be an arbitrary $n \times n$ permutation matrix. It is not difficult to check that the assumption of $P$ avoiding $A$ leads to $\mathrm{diag}(P^{-1}A) = (0, \ldots, 0)$. Thus by the former case we find $A_1', A_2'$ with

$$A_1' + A_2' = P^{-1}A, \text{ such that } A_1' + I_n, A_2' - I_n \in \mathrm{GL}_n(S).$$

Setting $A_1 := PA_1'$ and $A_2 := PA_2'$ completes the proof. $\square$

**Lemma 4.8** ([15, Proposition 4.4]). *Let $B = \mathrm{diag}(B_1, \ldots, B_t)$ be an $n \times n$ matrix, the size of each block $\leq n/2$. Then there is a permutation matrix $P$ that avoids $B$.*

*Proof.* Let $A$ be a matrix of size $\leq 3$, then the block size is at most one, indicating that $A$ is a diagonal matrix. Then

$$P := \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$

avoids $A$. Therefore we may assume that $n \geq 4$ and $t \geq 2$.

The idea of the proof is to develop an algorithm, which starting with an arbitrary permutation matrix $P$ outputs in every step a new permutation matrix $P'$, such that $m(P', B) < m(P, B)$. Thus eventually leading to a permutation matrix avoiding $B$. We arrange the proof into three steps.

*Preparation*:
Choose an arbitrary $n \times n$ permutation matrix $P$, suppose $m(P, B) > 0$. We find a pair of indices $(i, j)$ such that $P_{ij} \neq 0$ and $B_{ij} \neq 0$. A simultaneous multiplication of $P$ and $B$ by a permutation matrix does not alter the meeting number, hence we may assume that $(i, j)$ lies in the first block $B_1$ of size say $p \times q$. This relates to $i \leq p \leq n/2$ and $j \leq q \leq n/2$. Consider the following partitions, where the subindex of each block indicates its size:

$$B = \begin{pmatrix} B_{1\,p \times q} & \mathbf{0}_{p \times (n-q)} \\ \mathbf{0}_{(n-p) \times q} & A_{(n-p) \times (n-q)} \end{pmatrix}, P = \begin{pmatrix} P_{1\,p \times q} & P_{2\,p \times (n-q)} \\ P_{3\,(n-p) \times q} & P_{4\,(n-p) \times (n-q)} \end{pmatrix}$$

*Sudoku-condition*:
Regarding the above partition we prove that some entry of $P_4$ is 1: Suppose $P_4$ has only zero entries, then $P_3$ needs to contain $n - p$ non-zero elements, all fitting inside the first $q$ columns of $P$. In fact these $n - p$ elements even have to fit into the first $q - 1$ columns of $P$, as the $j$-th column is already reserved for the non-zero $(i, j)$-entry. As the numbers of columns and rows in $B_1$ is less or equal to $n/2$, we have $q - 1 < n/2$ and $p \leq n/2$ implying $q - 1 < n - p$. By the Pigeonhole principle there is a column within the first $q - 1$ columns featuring two non-zero elements - a contradiction to the definition of a permutation matrix.

*Decrement of meeting number*:
Now let $(r, c)$, with $r > p, c > q$, denote the position of the non-zero element of $P_4$, the existence of which we have just deduced. Let $E$ denote the permutation matrix that swaps rows $i$ and $r$ and put $P' = EP$. Now $P'$ has 1's in positions $(r, j)$ and $(i, c)$, but not at positions $(r, c)$ and $(i, j)$. Both positions $(r, j)$ and $(i, c)$ lie in the 0-blocks of $B$, hence $m(P', B) < m(P, B)$. □

*Proof of Proposition 4.4.* The claim follows readily from Lemma 4.7 and Lemma 4.8. □

## 4.2 Modules over Dedekind domains

### 4.2.1 Projective modules

For the proof of Theorem 4.3 we need to provide the theory of *finitely generated* and in particular *finitely presented* modules over Dedekind domains. Indeed these results will generalise the endeavours made in Chapter 2 regarding the structure theorem over principal ideal domains 2.25.

The primary key in obtaining these structure theorems is the fact, that a finitely generated module over a Dedekind domain is projective; this allows us to employ a broader technical framework based on the theory of projective modules, which we will now develop.

In this section let $R$ denote an arbitrary commutative ring.

**Definition 4.9** ([34, §4 Definition 3.1]). An $R$-module $P$ is called *projective*, if for arbitrary $R$-modules $A, B$ and module homomorphisms $f, g$ there exists a module

homomorphism $h$ turning the following diagram with exact row into a commutative diagram.

$$
\begin{array}{ccc}
 & & P \\
 & h \nearrow & \downarrow f \\
A & \xrightarrow{\;g\;} & B \longrightarrow 0
\end{array}
$$

**Proposition 4.10** ([34, §4 Proposition 3.5])**.** *Let* $\{P_i\}_{i \in I}$ *be a family of R-modules. Then* $\bigoplus_{i \in I} P_i$ *is projective, iff every* $P_i$ *is projective.*

*Proof.* "$\Rightarrow$" Take an arbitrary $P_0 \in \{P_i\}_{i \in I}$, let $g : A \to B$ be an epimorphism, $f_0 : P_0 \to B$ a module morphism. Furthermore denote by $\pi_0 : \bigoplus P_i \to P_0$ and $\iota_0 : P_0 \to \bigoplus P_i$ the canonical projection and injection respectively. The commutative diagram

$$
\begin{array}{ccc}
 & & \bigoplus P_i \\
 & \gamma \nearrow & \iota_0 \uparrow \downarrow \pi_0 \\
 & & P_0 \\
 & & \downarrow f_0 \\
A & \xrightarrow{\;g\;} & B
\end{array}
$$

shows the existence of $\gamma$ and that $\gamma \circ \iota_0$ may serve as the required homomorphism from $P_0$ to $A$.

"$\Leftarrow$" Let $f : \bigoplus P_i \to B$ and an epimorphism $g : A \to B$ be given. Consider the commutative diagram warranting homomorphisms $f_i$ for all $i \in I$ as $P_i$ is projective

$$
\begin{array}{ccc}
 & & P_i \\
 & f_i \nearrow & \downarrow \iota_i \\
 & & \bigoplus P_i \\
 & & \downarrow f \\
A & \xrightarrow{\;g\;} & B
\end{array}
$$

Then $(x_i)_{i \in I} \mapsto \sum f_i(x_i)$ from $\bigoplus P_i \to A$ is the required homomorphism.     $\square$

It is easy to see, that the notion of projective modules generalises free modules.

**Corollary 4.11** ([35, Corollary to Proposition 1.33]). *A free R-module $\mathcal{F}$ is projective.*

*Proof.* Since $\mathcal{F}$ is free, it is isomorphic to a sum of rings $R$. The last proposition shows, that it therefore suffices to prove that $R$ is projective. Consider

$$
\begin{array}{ccc}
 & R & \\
{\scriptstyle h} \swarrow & \big\downarrow {\scriptstyle f} & \\
A \xrightarrow{\ g\ } & B \longrightarrow & 0
\end{array}
$$

with arbitrary $R$-modules A,B and exact row. Suppose $f(1) = b$, select an $a \in A$ such that $g(a) = b$. Then $h : R \to A$ defined by $x \mapsto xa$ makes the diagram commutative, showing that $R$ is projective. $\qquad\square$

**Remark 4.12.** As we will be needing the universal property of free modules in the next proposition, we recall:

An $R$-module $\mathcal{F}$ is free with basis $X$, if for all $R$-modules $M$ and an arbitrary mapping $f : X \to M$ the latter can be lifted to a $R$-module morphism $\overline{f} : \mathcal{F} \to M$, making the following diagram commutative:

$$
\begin{array}{ccc}
\mathcal{F} & & \\
\big\uparrow\big\downarrow & \searrow^{\overline{f}} & \\
X & \xrightarrow{\ f\ } & M
\end{array}
\tag{4.1}
$$

**Proposition 4.13** ([34, §4 Proposition 3.4], [35, Proposition 1.34 and 1.35]). *Let P be an R-module. The following are equivalent:*

   (i) *P is projective.*

   (ii) *If $0 \to A \xrightarrow{i} B \xrightarrow{p} P \to 0$ is a short exact sequence, then $B \cong A \oplus P$.*

   (iii) *There exists a free module $\mathcal{F}$ and an R-module K such that $\mathcal{F} \cong K \oplus P$.*

   (iv) *There exist elements $\{a_t\}_{t \in T} \subseteq P$ and homomorphisms $\{f_t\}_{t \in T} \subseteq \mathrm{Hom}_R(P, R)$, such that every $a \in P$ may be written as*

$$
a = \sum_{t \in T} f_t(a) a_t,
$$

   *where only finitely many $f_t(a) \neq 0$.*

*Proof.*

$(i) \Rightarrow (ii)$. Since $P$ is projective, there exists a monomorphism $f : P \to B$ turning

$$
\begin{array}{ccc}
 & & P \\
 & \overset{f}{\nearrow} & \downarrow \text{id} \\
0 \longrightarrow A \overset{i}{\longrightarrow} B \overset{p}{\longrightarrow} P \longrightarrow 0
\end{array}
$$

into a commutative diagram.[2] As $i(A) \cong A$ and $f(P) \cong A$ it suffices to show that $B = i(A) \oplus f(P)$. Let $b \in B$ be arbitrary, then

$$p\big(b - (f \circ p)(b)\big) = p(b) - (\text{id} \circ p)(b) = 0$$

and hence $b - (f \circ p)(b) \in \ker(p)$ leading to

$$b = \big(b - (f \circ p)(b)\big) + (f \circ p)(b) \in \ker(p) + \text{im}(f) \subseteq B.$$

To see that directness of the sum take $b \in \ker(p) + \text{im}(f)$, then $p(b) = 0$ and $\exists a \in P : f(a) = b$. This leads to $0 = p(b) = (p \circ f)(a) = a$, finishing the proof.

$(ii) \Rightarrow (iii)$. Let $\mathcal{F}$ be a free $R$-module with basis $P$, such that the universal property (4.1) leads to the commutative diagram

$$
\begin{array}{ccc}
 & \mathcal{F} & \\
\iota \uparrow & \overset{f}{\searrow} & \\
P & \overset{\text{id}}{\longrightarrow} & P
\end{array}
$$

We derive an exact sequence

$$0 \to \ker(f) \to \mathcal{F} \to P \to 0.$$

By $(ii)$ we see $\mathcal{F} \cong \ker(f) \oplus P$.

$(iii) \Rightarrow (iv)$. Keeping the notation from the last step, let $\{x_t\}_{t \in T} \subseteq \mathcal{F}$ be a system of generators of $\mathcal{F}$. For $a \in P$, we have $\iota(a) = \sum_{t \in T} f_t(a) x_t$ with suitable $f_t(a) \in R$, the right hand side containing only finitely many non-zero terms. Applying $f$ yields $a = \sum_{t \in T} f_t(a) f(x_t)$ for all $a \in P$. As the $f_t : P \to P$ are homomorphisms due to their definition, setting $a_t = f(x_t)$ verifies the claim.

$(iv) \Rightarrow (iii)$. Denote by $\{x_t\}_{t \in T}$ the generators of some free module $\mathcal{F}$. Define a homomorphism $f : \mathcal{F} \to P$ by $f(x_t) = a_t$. Let the monomorphism $g : P \to \mathcal{F}$ be given by $a = \sum_{t \in T} f_t(a) a_t \mapsto \sum_{t \in T} f_t(a) x_t$. For any $a \in P$ we have

$$(f \circ g)(a) = f\Big( \sum_{t \in T} f_t(a) x_t \Big) = \sum_{t \in T} f_t(a) a_t = a,$$

---

[2] In terms of homological algebra $f$ constitutes a right split of the short exact sequence.

indicating that the composition $P \xrightarrow{g} \mathcal{F} \xrightarrow{f} P$ equals the identity map. Whence $P$ is a direct summand of $\mathcal{F}$.

$(iii) \Rightarrow (i)$. Suppose $K \oplus P \cong \mathcal{F}$, where $\mathcal{F}$ is free and a fortiori projective by Corollary 4.11. Proposition 4.10 shows that $P$ (and $K$) must be projective. $\qquad \square$

**Example 4.14.** Examples of projective modules may be found in any textbook about abstract algebra, we give a few outlines

- Clearly every vector space is free and hence projective.
- Set $R = \mathbb{Z}/6\mathbb{Z}$, then $R$ is a free module over itself. As $R \cong \mathbb{Z}/2\mathbb{Z} \oplus \mathbb{Z}/3\mathbb{Z}$, we deduce from Proposition 4.13(ii), that $\mathbb{Z}/2\mathbb{Z}$ is projective over $R$. Moreover by simply comparing the number of elements $\mathbb{Z}/2\mathbb{Z}$ can not be a free $R$-module. Note that invoking the Chinese Remainder Theorem on more general rings warrants a plethora of similar examples.

Recall that an element $x \in P$ is said to be *torsion*, if $\exists r \in R \setminus \{0\} : rx = 0$.

**Lemma 4.15** ([35, Lemma 1.37 and 1.38]). *Let $\mathcal{D}$ be an integral domain such that each of its ideals is projective. Suppose $P$ is a finitely generated, torsion-free $\mathcal{D}$-module, then it is isomorphic to a finite direct sum of ideals of $\mathcal{D}$.*

*Proof.* Denote by $a_1, \ldots, a_m$ a set of generators of $P$, and let $K$ be the field of fractions of $\mathcal{D}$. Let $y_1, \ldots, y_n$ be a basis of the finite dimensional $K$-vector space $Ka_1 + \cdots + Ka_m$. There exist $r_{ij} \in K$ such that $a_i = \sum_{j=1}^{n} r_{ij} y_j$. We find an element $q \in \mathcal{D}$ that clears out the denominators, i.e $qr_{ij} \in \mathcal{D}$. In view of

$$a_i = \sum_{j=1}^{n} (qr_{ij}) \frac{y_j}{q} \in \mathcal{D}\frac{y_1}{q} + \cdots + \mathcal{D}\frac{y_n}{q} := \mathcal{F}$$

we arrive at

$$P = \mathcal{D}a_1 + \cdots \mathcal{D}a_m \subseteq \mathcal{F},$$

where $\mathcal{F}$ constitutes a free $\mathcal{D}$-module. The rank of $\mathcal{F}$ equals $n$, as the $\frac{y_1}{q}, \ldots, \frac{y_n}{q}$ are $K$-linearly independent and thus also $\mathcal{D}$-linearly independent.
In the case of $\mathcal{F}$ featuring rank 0, there is nothing to prove, since $P$ must be 0. Suppose thus the claim holds for all $\mathcal{D}$-modules contained in some free $\mathcal{D}$-module of rank $n-1$. Let $x_1, \ldots, x_n$ denote the free generators of $\mathcal{F} \supseteq P$. Denote by $\mathcal{F}_{n-1}$ the free $\mathcal{D}$-module generated by $x_1, \ldots, x_{n-1}$. Let $P \ni x = \sum_{i=1}^{n} r_i x_i$ with $r_i \in \mathcal{D}$. Then $f : x \mapsto r_n \in \operatorname{Hom}_{\mathcal{D}}(P, \mathcal{D})$ is well-defined as $P$ is torsion-free. Due to assumption the ideal $\operatorname{im}(f)$ of $\mathcal{D}$ in the exact sequence

$$0 \to \ker(f) \to P \to \operatorname{im}(f) \to 0$$

is projective, hence by Proposition 4.10 we see $P \cong \mathrm{im}(f) \oplus \ker(f)$. As $\ker(f)$ is contained in $\mathcal{F}_{n-1}$, the induction hypothesis is applicable, showing that $P$ is a direct sum of ideals in $\mathcal{D}$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\square$

### 4.2.2 Structure theorems for modules over Dedekind domains

In this section $\mathfrak{O}$ denotes a Dedekind domain and $K$ its field of fractions.

**Proposition 4.16** ([35, Lemma 1.36]). *A non-zero ideal $\mathfrak{a}$ in an integral domain $\mathcal{D}$ is projective, iff it is invertible.*

*Proof.* Let $K$ denote the field of fractions of $\mathcal{D}$.
"$\Rightarrow$" Assume $\mathfrak{a}$ is projective, let $\{a_t\}_{t \in T} \subseteq \mathfrak{a}$ and let $\{f_t\}_{t \in T}$ be homomorphisms as in Proposition 4.13(iv). Take an arbitrary $x \in \mathfrak{a}$, one has $f_t(1)x = f_t(x) \in \mathcal{D}$, since $f_t$ is $\mathcal{D}$-linear. Therefore $f_t(1)\mathfrak{a} \subseteq \mathcal{D}$ and thus

$$f_t(1) \in \mathfrak{a}' := \{\alpha \in K | \alpha\mathfrak{a} \subseteq \mathcal{D}\}.$$

We compute

$$x = \sum_{t \in T} f_t(x)a_t = x \sum_{t=1}^{n} f_t(1)a_t,$$

so that $1 = \sum_{t=1}^{n} f_t(1)a_t$. The transition from a sum over an infinite index set to a finite sum is justified due to proposition 4.13(iv). This shows $\mathcal{D} \subseteq \mathfrak{a}\mathfrak{a}'$ and hence $\mathfrak{a}$ is invertible.

"$\Leftarrow$" Suppose $\mathfrak{a}\mathfrak{a}^{-1} = D$, then we find $a_1, \dots, a_n \in \mathfrak{a}$ and $x_1, \dots, x_n \in \mathfrak{a}^{-1}$ such that $\sum_{t=1}^{n} x_i a_i = 1$. Take an arbitrary $x \in \mathfrak{a}$, then

$$x = \sum_{t=1}^{n} x x_t a_t = \sum_{t=1}^{n} f_t(x)a_t,$$

where we set $f_t(x) = x x_t$. Now the elements $a_t$ and homomorphisms $f_t$ are as required by Proposition 4.13(iv), hence $\mathfrak{a}$ is projective. $\qquad\qquad\qquad\qquad\qquad\square$

The consequent corollary follows directly from the proposition just proved; it opens the door for treating ideal theoretic problems in Dedekind domains via the theory of projective modules.

**Corollary 4.17.** *All ideals in a Dedekind domain are projective.*

A commutative ring $R$ is called *(semi-)hereditary*, if all its (finitely generated) modules are projective over $R$. The proposition shows that the notions of Dedekind domain and hereditary integral domain coincide, thus exhibiting another characterisation of

Dedekind domains via projective modules. Moreover a Prüfer domain is the same as a semihereditary integral domain.

**Definition 4.18.** Given a finitely generated $\mathfrak{O}$-module $M$ over a Dedekind domain $\mathfrak{O}$, we write $\mathcal{T}(M)$ to denote the submodule of $M$ containing all torsion elements of $M$.

**Theorem 4.19** ([35, Theorem 1.32]). *For a finitely generated $\mathfrak{O}$-module $M$, we find $k \in \mathbb{N}$ and an ideal $\mathfrak{d}$ such that $M \cong \mathfrak{O}^k \oplus \mathfrak{d} \oplus \mathcal{T}(M)$.*

*Proof.* Defining $M_1 = {}^{M}\!/_{\mathcal{T}(M)}$, we see that $M_1$ is a torsion-free and finitely generated $\mathfrak{O}$-module. As all modules over $\mathfrak{O}$ are projective by Corollary 4.17, we invoke Lemma 4.15 to write $M_1 = \mathfrak{d}_1, \ldots, \mathfrak{d}_{k+1}$ as a sum of ideals $\mathfrak{d}_i$, i.e projective modules, which entails by Proposition 4.10, that $M_1$ is projective over $\mathfrak{O}$ itself. As

$$0 \to \mathcal{T}(M) \to M \to M_1 \to 0$$

is exact, Proposition 4.13 shows that

$$M \cong \mathcal{T}(M) + M_1 \cong \mathcal{T}(M) \oplus \bigoplus_{j=1}^{k+1} \mathfrak{d}_j.$$

The proof is complete, if we derive the existence of an ideal $\mathfrak{d}$, such that $\bigoplus_{j=1}^{k+1} \mathfrak{d}_j \cong \mathfrak{O}^k \oplus \mathfrak{d}$; this fact will be a special case of the next theorem. $\qquad\square$

The former theorem is due to Steinitz [5] who proved it in 1912 for rings of algebraic integers. The upcoming theorem, which will play a crucial role in the proof of the structure theorem for finitely presented modules as well, is eo ipso of high interest as it provides an alternative method to handle direct sums of ideals in Dedekind domains.

We slide in the auxiliary

**Lemma 4.20** ([35, cf. Corollary 6 to Proposition 1.14]). *For any two ideals $\mathfrak{a}, \mathfrak{b}$ of a Dedekind domain $\mathfrak{O}$, we find an ideal $\mathfrak{a}'$ of $\mathfrak{O}$ satisfying $[\mathfrak{a}'] = [\mathfrak{a}]$, such that $(\mathfrak{a}', \mathfrak{b}) = \mathfrak{O}$.*

*Proof.* Let $\mathfrak{a} = \prod \mathfrak{p}^{\alpha_{\mathfrak{p}}}$ be the factorisation of $\mathfrak{a}$ into prime ideals. Denote by $\mathfrak{r}_1, \ldots, \mathfrak{r}_n$ the prime ideals dividing $\mathfrak{b}$ and not dividing $\mathfrak{a}$. Choose non-zero elements $x_{\mathfrak{p}} \in \mathfrak{p}^{\alpha_{\mathfrak{p}}} \setminus \mathfrak{p}^{\alpha_{\mathfrak{p}}+1}$ for all $\mathfrak{p}$ dividing $\mathfrak{a}$. Finally select a prime ideal $\mathfrak{p}_0$ not dividing $\mathfrak{a}$ or $\mathfrak{b}$. The Chinese remainder theorem yields a solution $u$ to the system of congruences with

pairwise relatively prime moduli:

$$u \equiv x_{\mathfrak{p}} \pmod{\mathfrak{p}^{\alpha_{\mathfrak{p}}+1}} \quad \forall \mathfrak{p} : \mathfrak{p} | \mathfrak{a}$$
$$u \equiv 1 \pmod{\mathfrak{r}_i} \qquad \forall 1 \leq i \leq n$$
$$u \equiv 0 \pmod{\mathfrak{p}_0}.$$

Clearly $u \notin \mathfrak{p}^{\alpha_{\mathfrak{p}}+1}$, but there exists $x'_{\mathfrak{p}} \in \mathfrak{p}^{\alpha_{\mathfrak{p}}+1}$ such that $u = x_{\mathfrak{p}} + x'_{\mathfrak{p}} \in \mathfrak{p}^{\alpha_{\mathfrak{p}}}$. Thus $\mathfrak{p}^{\alpha_{\mathfrak{p}}} \| u\mathfrak{O}$ for each $\mathfrak{p}$ dividing $\mathfrak{a}$, which leads to $\mathfrak{a} | \mathfrak{O}$. We find an ideal $\mathfrak{c}$ such that

$$u\mathfrak{O} = \mathfrak{c} \prod \mathfrak{p}^{\alpha_{\mathfrak{p}}} = \mathfrak{c}\mathfrak{a}.$$

We observe that $\mathfrak{p} \nmid \mathfrak{c}$, as $\mathfrak{p}^{\alpha_{\mathfrak{p}}}$ is the highest power of $\mathfrak{p}$ in $u\mathfrak{O}$. Moreover for $1 \leq i \leq n$ we have $\mathfrak{r}_i \nmid \mathfrak{c}$ as $\mathfrak{r}_i \nmid u\mathfrak{O}$ and additionally $\mathfrak{c} \neq \mathfrak{O}$ since $\mathfrak{p}_0 | \mathfrak{c}$. This shows that $\mathfrak{c}$ and $\mathfrak{b}$ are relatively prime.

Let $\mathfrak{c} = \prod_{i=1}^{t} \mathfrak{q}_i^{\beta_i}$ be the prime factorisation of $\mathfrak{c}$. As $\mathfrak{c} + \mathfrak{b} = \mathfrak{O}$ invoking again the Chinese remainder theorem, we obtain $b \in \mathfrak{q}_i^{\beta_i} \setminus \mathfrak{q}_i^{\beta_i+1}$ $(1 \leq i \leq t)$ satisfying $b \equiv 1 \pmod{\mathfrak{b}}$. Whence $\mathfrak{c} | b\mathfrak{O}$, i.e. there is an ideal $\mathfrak{a}'$ such that $\mathfrak{c}\mathfrak{a}' = b\mathfrak{O}$ and moreover $1 = b + (1 - b) \in b\mathfrak{O} + \mathfrak{b}$. Thus

$$\mathfrak{O} = (b\mathfrak{O}, \mathfrak{b}) = (\mathfrak{c}\mathfrak{a}', \mathfrak{b}) = (\mathfrak{a}', \mathfrak{b})$$

and

$$[\mathfrak{a}] = [b\mathfrak{a}] = [(\mathfrak{c}\mathfrak{a}')\mathfrak{a}] = [(\mathfrak{c}\mathfrak{a})\mathfrak{a}'] = [u\mathfrak{a}'] = [\mathfrak{a}'].$$

$\square$

**Theorem 4.21** ([35, Theorem 1.39])**.** *Let $M_1, M_2$ be torsion-free $\mathfrak{O}$-modules, given by $M_1 = \bigoplus_{\ell=1}^{m} \mathfrak{a}_\ell$ and $M_2 = \bigoplus_{\ell=1}^{n} \mathfrak{b}_\ell$, where $\mathfrak{a}_\ell, \mathfrak{b}_\ell$ are ideals of $\mathfrak{O}$. Then $M_1$ and $M_2$ are isomorphic as $\mathfrak{O}$-modules if and only if*

$$n = m \text{ and } \Big[ \prod_{\ell=1}^{m} \mathfrak{a}_\ell \Big] = \Big[ \prod_{\ell=1}^{n} \mathfrak{b}_\ell \Big].$$

*Proof.*
"$\Leftarrow$" It suffices to evidence the claim for $m = n = 2$, as the other cases may be settled via induction. Lemma 4.20 shows that, we find an ideal $\mathfrak{a}'_1$ with $[\mathfrak{a}'_1] = [\mathfrak{a}_1]$ such that $\mathfrak{a}'_1$ and $\mathfrak{a}_2$ are relatively prime. As $\mathfrak{O}$ is projective, the following commutative diagram

$$
\begin{array}{ccccccccc}
0 & \longrightarrow & \mathfrak{a}'_1 \cap \mathfrak{a}_2 & \longrightarrow & \mathfrak{a}'_1 \oplus \mathfrak{a}_2 & \longrightarrow & \mathfrak{a}'_1 + \mathfrak{a}_2 & \longrightarrow & 0 \\
 & & \Big\| & & \Big| & & \Big\| & & \\
0 & \longrightarrow & \mathfrak{a}'_1 \mathfrak{a}_2 & \longrightarrow & \mathfrak{a}'_1 \oplus \mathfrak{a}_2 & \longrightarrow & \mathfrak{O} & \longrightarrow & 0
\end{array}
$$

with exact rows warrants

$$[\mathfrak{a}_1 \oplus \mathfrak{a}_2] = [\mathfrak{a}_1' \oplus \mathfrak{a}_2] = [\mathfrak{O} \oplus \mathfrak{a}_1' \mathfrak{a}_2] = [\mathfrak{O}] + [\mathfrak{a}_1' \mathfrak{a}_2].$$

The same process yields an ideal $\mathfrak{b}_1'$ with $[\mathfrak{b}_1'] = [\mathfrak{b}_1]$ such that

$$[\mathfrak{b}_1 \oplus \mathfrak{b}_2] = [\mathfrak{b}_1' \oplus \mathfrak{b}_2] = [\mathfrak{O}] + [\mathfrak{b}_1' \mathfrak{b}_2].$$

From the assumption $[\mathfrak{a}_1 \mathfrak{a}_2] = [\mathfrak{b}_1 \mathfrak{b}_2]$ we immediately get $[\mathfrak{a}_1 \oplus \mathfrak{a}_2] = [\mathfrak{b}_1 \oplus \mathfrak{b}_2]$. Thus there exists an $0 \neq x \in K$ such that $\mathfrak{a}_1 \oplus \mathfrak{a}_2 = x(\mathfrak{b}_1 \oplus \mathfrak{b}_2)$, implying $M_1 \cong M_2$.

"$\Rightarrow$" In the commutative diagram

$$\operatorname{span}_K \big(\iota_1(M_1)\big) = K^m \xleftarrow{\overline{\varphi} \;\cong\;} K^n = \operatorname{span}_K \big(\iota_2(M_2)\big)$$



the embeddings $\iota_1, \iota_2$ induced by the canonical embeddings depicted in the diagram as outer arrows, enables the lift of $\varphi$ to a $K$-morphism $\overline{\varphi}$. A comparison of dimensions yields $n = m$.

Let us assume that $\mathfrak{O} \subseteq \mathfrak{a}_i, \mathfrak{b}_i$ for $1 \leq i \leq m$, this condition poses no restriction: for an integral ideal $\mathfrak{a}_i$ we find a non-zero $x \in K$, such that $\mathfrak{O} \subseteq x\mathfrak{a}_i := \mathfrak{a}_i'$, which entails $[\mathfrak{a}_i] = [\mathfrak{a}_i']$.

Now denote again by $\varphi$ the isomorphism $M_1 \to M_2$ and let $\varphi_i = \varphi_{|\mathfrak{a}_i}$ the restriction of $\varphi$ to $\mathfrak{a}_i$. If $a, x \in \mathfrak{a}_i$, we have linearity: $\varphi_i(xa) = x\varphi_i(a)$. Indeed let $x = \frac{\alpha}{\beta}$ with $\alpha, \beta \in \mathfrak{O}$, then

$$\beta\varphi_i(ax) = \beta\varphi_i\Big(a\frac{\alpha}{\beta}\Big) = \alpha\varphi_i(a).$$

Let $(a_{i1}, \ldots, a_{im})$ be the image of $1$ under $\varphi_i$, where $a_{ij} \in \mathfrak{b}_i$ for $1 \leq i, j \leq m$. Moreover let $\pi_j$ be the projection from $M_2$ onto $\mathfrak{b}_j$, we compute

$$\sum_{i=1}^m a_{ij}\mathfrak{a}_i = \Big\{ \sum_{i=1}^m a_{ij}x_i \,|\, x_i \in \mathfrak{a}_i \Big\} = \Big\{ \pi_j\Big( \sum_{i=1}^m a_{i1}x_i, \ldots, \sum_{i=1}^m a_{im}x_i \Big) \,|\, x_i \in \mathfrak{a}_i \Big\}$$

$$= \Big\{ \pi_j\Big( \sum_{i=1}^m (a_{i1}, \ldots, a_{im})x_i \Big) \,|\, x_i \in \mathfrak{a}_i \Big\} = \Big\{ \pi_j\Big( \sum_{i=1}^m \underbrace{\varphi_i(1)x_i}_{\varphi_i(x_i)} \Big) \,|\, x_i \in \mathfrak{a}_i \Big\}$$

$$= \Big\{ (\pi_j \circ \varphi)\Big( \sum_{i=1}^m x_i \Big) \,|\, x_i \in \mathfrak{a}_i \Big\} = \mathfrak{b}_j \quad \text{for } 1 \leq j \leq m.$$

For some permutation on $m$ letters $\sigma \in \mathfrak{S}_m$ put $C_\sigma = \prod_{i=1}^m a_{\sigma(i),i}$. In view of

$$\prod_{j=1}^m \mathfrak{b}_j = \prod_{j=1}^m \sum_{i=1}^m a_{ij} \mathfrak{a}_i$$

combinatorial considerations lead to the insight that the coefficient of $\prod_{i=1}^m \mathfrak{a}_i$ is equal to $\sum_{\sigma \in \mathfrak{S}_m} C_\sigma$. Therefore

$$\sum_{\sigma \in \mathfrak{S}_m} C_\sigma \prod_{i=1}^m \mathfrak{a}_i \subseteq \prod_{j=1}^m \mathfrak{b}_j.$$

In particular, for any $\sigma \in \mathfrak{S}_m$ and $x_i \in \mathfrak{a}_i$ we have

$$\mathrm{sgn}(\sigma) \cdot C_\sigma \prod_{i=1}^m x_i \in \prod_{j=1}^m \mathfrak{b}_j.$$

Finally

$$\sum_{\sigma \in \mathfrak{S}_m} \prod_{i=1}^m \mathrm{sgn}(\sigma) C_\sigma x_i,$$

which due to the Leibniz determinant formula equals $\det(a_{ij}) \prod_{i=1}^m x_i$, is contained in $\prod_{j=1}^m \mathfrak{b}_j$.

Repeating the process with $M_1$ and $M_2$ interchanged, we find a matrix $(b_{ij})_{1 \le i,j \le m}$ granting

$$\det(b_{ij}) \prod_{j=1}^m \mathfrak{b}_j \subseteq \prod_{i=1}^m \mathfrak{a}_i.$$

As in fact these matrices are inverse to each other, we arrive at

$$\prod_{i=1}^m \mathfrak{a}_i = \det(b_{ij}) \prod_{j=1}^m \mathfrak{b}_j.$$

$\square$

**Remark 4.22.** Employing the theorem just established, and broader ideal- and module-theoretic results [35, §1.3.3] the torsion module in Theorem 4.19 may be further decomposed, such that for a finitely generated module $M$ over $\mathfrak{O}$, we have

$$M \cong \bigoplus_{j=1}^s \mathfrak{O}\big/_{\mathfrak{a}_j} \oplus \mathfrak{d} \oplus \mathfrak{O}^k,$$

for suitable integral ideals $\mathfrak{a}_1, \ldots, \mathfrak{a}_s$ and some fractional ideal $\mathfrak{d}$ of $\mathfrak{O}$.

## 4.3 Indecomposable matrices over Dedekind domains

This section is devoted to the proof of Theorem 4.3, which bounds the size of an indecomposable matrix by the class number. We begin by introducing finitely presented modules over Dedekind domains $\mathfrak{D}$, which are a special case of finitely generated modules.

**Definition 4.23** ([4, Definition in §1])**.** Let $A$ be an $m \times n$ matrix over a commutative ring $\mathfrak{D}$. $A$ induces an $\mathfrak{D}$-module homomorphism $\varphi : \mathfrak{D}^m \to \mathfrak{D}^n$ via $x \mapsto xA$. Denoting by $M_A$ the image of $\varphi$ in $\mathfrak{D}^n$, i.e. $M_A$ is generated by the not necessarily linearly independent rows of $A$, we obtain an exact sequence

$$\mathfrak{D}^m \xrightarrow{\varphi} \mathfrak{D}^n \twoheadrightarrow S_A := \mathfrak{D}^n\!/\!M_A \to 0.$$

A module $S_A$ arising in this way is called *finitely presented*.

Levy [4] utilised the work of Krull [6] to derive a structure theorem for finitely presented modules.

**Theorem 4.24** ([4, Separated Divisor Theorem])**.** *Let $A$ be an $m \times n$ matrix of rank $r$ over $\mathfrak{D}$, then $S_A$ takes the form*

$$S_A = S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{d})_{(m \times n)} := \begin{cases} \bigoplus_{j=1}^r \mathfrak{D}\!/\!\mathfrak{h}_j \oplus \mathfrak{d} \oplus \mathfrak{D}^{n-r-1} & \text{if } r < n \\ \bigoplus_{j=1}^r \mathfrak{D}\!/\!\mathfrak{h}_j & \text{if } r = n, \end{cases}$$

*where the $\mathfrak{h}_j$'s are integral ideals and $\mathfrak{d}$ is a fractional ideal of $\mathfrak{D}$. Furthermore*

$$\left[ \prod_{j=1}^r \mathfrak{h}_j \right] = [\mathfrak{d}] \qquad\qquad \text{if } r = m \qquad\qquad \text{(a)}$$

$$[\mathfrak{d}] = [\mathfrak{D}] \qquad\qquad \text{if } r = n \text{ or } r = 0 \qquad\qquad \text{(b)}$$

*the latter being just a notational convention. Furthermore for given integral ideals $\mathfrak{h}_j$, fractional ideal $\mathfrak{d}$ and positive integers $m, n$ there exists always a matrix $A$ satisfying $S_A = S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{d})_{(m \times n)}$.*

We take a look at a simple example.

**Example 4.25.** Setting $\mathfrak{D} = \mathbb{Z}$, $A = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}$ and $M_A = \operatorname{im}(\varphi) \in \mathbb{Z}^2$, we derive

$$S_A \cong \mathbb{Z}^2\!/\!M_A \cong \mathbb{Z}^2\!/\!(2,0)\mathbb{Z} + (0,3)\mathbb{Z} \cong \mathbb{Z}\!/\!2\mathbb{Z} \oplus \mathbb{Z}\!/\!3\mathbb{Z} \cong \mathbb{Z}\!/\!6\mathbb{Z}.$$

This calculation shows that the decomposition obtained via the Separated Divisor Theorem 4.24 is in general not unique.

In order to utilise the Separated Divisor Theorem to evidence Theorem 4.3, we need the following result.

**Theorem 4.26** ([36])**.** *For a Dedekind domain $\mathfrak{O}$ let $A, B \in \mathrm{Mat}_n(\mathfrak{O})$, then $A \sim B \Leftrightarrow S_A \cong S_B$.*

The result is essentially due to Levy, who proved it first by taking advantage of Krull's work [6]. In the follow-up paper [36] Levy and Robson generalised the statement by using extensive module-theoretic methods.

**Lemma 4.27** (Diagonalisation Lemma [4])**.** *Let $A$ be an $m \times n$ matrix of rank $r$ and $B$ a $p \times q$ matrix of rank $s$ over $\mathfrak{O}$. In view of Theorem 4.24 we write*

$$S_A = S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{o}) \text{ and } S_B = S(\mathfrak{h}'_1, \ldots, \mathfrak{h}'_s; \mathfrak{o}').$$

*Then*
$$S_{\mathrm{diag}(A,B)} \cong S_A \oplus S_B \cong S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r, \mathfrak{h}'_1, \ldots, \mathfrak{h}'_s; \mathfrak{o}\mathfrak{o}')_{(m+p) \times (n+q)}.$$

*Proof.* Similarly to the proof of Theorem 2.25 we find $M_{\mathrm{diag}(A,B)} = M_A \oplus M_B$, as these modules are generated by the rows of $A$ and $B$ respectively. We compute

$$S_{\mathrm{diag}(A,B)} \cong \mathfrak{O}^{n+q} \big/ M_{\mathrm{diag}(A,B)} \cong \mathfrak{O}^n \big/ M_A \oplus \mathfrak{O}^q \big/ M_B \cong S_A \oplus S_B.$$

To evidence the isomorphy for the second term in the claim, we calculate

$$S_A \oplus S_B \cong S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{o})_{(m \times n)} \oplus S(\mathfrak{h}'_1, \ldots, \mathfrak{h}'_s; \mathfrak{o}')_{(p \times q)}$$
$$= \bigoplus_{i=1}^{r} \mathfrak{O} \big/ \mathfrak{h}_i \oplus \bigoplus_{i=1}^{s} \mathfrak{O} \big/ \mathfrak{h}'_i \oplus \mathfrak{o} \oplus \mathfrak{o}' \oplus \mathfrak{O}^{n-r-1} \oplus \mathfrak{O}^{q-s-1}.$$

Now Theorem 4.21 shows that $\mathfrak{o} \oplus \mathfrak{o}' \cong \mathfrak{O} \oplus \mathfrak{o}\mathfrak{o}'$, hence we arrive at

$$S_{\mathrm{diag}(A,B)} \cong S_A \oplus S_B$$
$$\cong \bigoplus_{i=1}^{r} \mathfrak{O} \big/ \mathfrak{h}_i \oplus \bigoplus_{i=1}^{s} \mathfrak{O} \big/ \mathfrak{h}'_i \oplus \mathfrak{o}\mathfrak{o}' \oplus \mathfrak{O}^{n+q-(r+s)-1}$$
$$= S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r, \mathfrak{h}'_1, \ldots, \mathfrak{h}'_s; \mathfrak{o}\mathfrak{o}')_{(m+p) \times (n+q)}.$$

$\square$

**Lemma 4.28.**

   *(i) Let $S_A$ be a finitely presented module induced by a matrix $A$ over $\mathfrak{O}$. If $S_A$ may be decomposed into a sum of two non-vanishing finitely presented modules, then $A$ is decomposable.*

   *(ii) Any $m \times n$ matrix $A$ over $\mathfrak{O}$ inducing a torsion module $S_A$ and satisfying $m > h_\mathfrak{O}$ or $n > h_\mathfrak{O}$ is decomposable.*

*Proof.*

(i) $S_A$ may be decomposed into a sum of two finitely presented modules, none of them vanishing. Due to the Separated Divisor Theorem 4.24 these two modules take the form $S_B$ and $S_C$ for suitable matrices $B, C$. Thus by the Diagonalisation Lemma 4.27 we find $S_A = S_B \oplus S_C = S_{\mathrm{diag}(B,C)}$. An application of Theorem 4.26 evidences $A \sim \mathrm{diag}(B,C)$.

(ii) Let $r$ be the rank of $A$. The Separated Divisor Theorem shows that $n = r$ in the case of $S_A$ being a torsion module.

- If $n = r > 1$, then

$$S_A = \bigoplus_{i=1}^{r} \mathfrak{O}\big/_{\mathfrak{h}_i} = \bigoplus_{i=1}^{r} S(\mathfrak{h}_i; \mathfrak{O})_{2\times 1}.$$

Invoking part one of the lemma verifies the claim.

- If $n = r = 1$, then $S_A$ has the form $S_A = \mathfrak{O}\big/_{M_A} = \mathfrak{O}\big/_{\mathfrak{h}}$ with some integral ideal $\mathfrak{h}$.

Now $m = 1$ is impossible due to $m > h_{\mathfrak{O}} \geq 1$ or $n > h_{\mathfrak{O}} \geq 1$

If $m = 2$, then $A = (a, b)^\mathsf{T}$ with $a, b \in \mathfrak{O}$. As $m > h_{\mathfrak{O}}$ or $n > h_{\mathfrak{O}}$ we must have $h_{\mathfrak{O}} = 1$. Thus $\mathfrak{O}$ constitutes a principal ideal domain and we may use the Smith Normal Form 2.22 to obtain $d \in \mathfrak{O}$, such that $A \sim (d, 0)^\mathsf{T} \sim \mathrm{diag}(d, \mathbf{0}_{1\times 0})$. This shows that $A$ is decomposable.

For $m \geq 2$ recall that every ideal in a Dedekind domain is generated by two elements. Thus $\mathfrak{O}\big/_{\mathfrak{h}} \cong \mathfrak{O}\big/_{x\mathfrak{O} + y\mathfrak{O}}$ for some elements $x, y \in \mathfrak{O}$. This shows that

$$A = (x, y, \underbrace{0, \ldots, 0}_{(m-2)\text{-times}})^\mathsf{T}.$$

Thus $A \sim \mathrm{diag}(x, y, \mathbf{0}_{(m-2)\times 0})$.

$\square$

**Lemma 4.29.** *Let $A$ be a $m \times n$ matrix of rank $r$ over $\mathfrak{O}$, satisfying $m \geq r + 2$ or $n \geq r + 2$. Moreover let $m > h_{\mathfrak{O}}$ or $n > h_{\mathfrak{O}}$. Then $A$ is decomposable, i.e. there exist two matrices $B, C$ over $\mathfrak{O}$, such that $A = \mathrm{diag}(B, C)$.*

*Proof.* Due to the prior Lemma 4.28(i) it suffices to show that $S_A$ is decomposable into a sum of two non-vanishing finitely presented modules. We need to consider some cases separately and will treat them in the following order.

| Case label | Condition 1 | Condition 2 |
|:---:|:---:|:---:|
| (i) | $n \geq r + 2$ | $m \geq r + 1$ |
| (ii) | $n \geq r + 2$ | $m = r$ |
| (iii) | $m \geq r + 2$ | $n \geq r + 2$ |
| (iv) | $m \geq r + 2$ | $n = r + 1$ |
| (v) | $m \geq r + 2$ | $n = r$ |

Note that we merely stated (iii) for convenience, as it is in fact covered by case (i).

(i) If $n \geq r + 2$ and $m \geq r + 1$, set $p = q = r + 1$. By means of the Separated Divisor Theorem 4.24 we calculate

$$S_A = S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{d})_{m \times n} = \bigoplus_{i=1}^{r} \mathfrak{D}/_{\mathfrak{h}_i} \oplus \mathfrak{d} \oplus \mathfrak{D}^{n-r-1}$$

$$= \bigoplus_{i=1}^{r} \mathfrak{D}/_{\mathfrak{h}_i} \oplus \mathfrak{d} \oplus \mathfrak{D}^{q-r-1} \oplus \mathfrak{D} \oplus \mathfrak{D}^{n-q-1}$$

$$= S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{d})_{p \times q} \oplus S(\varnothing; \mathfrak{D})_{(m-p) \times (n-q)}.$$

None of the "exceptional cases" in the Separated Divisor Theorem arise and the second summand does not vanish since $n - q = n - r - 1 \geq 1$.

(ii) When $n \geq r + 2$ and $m = r$, the previous decomposition works, if we put $p = r$ and $q = r + 1$: The condition (a) of the Separated Divisor Theorem is fulfilled for the first summand. The second summand does not vanish since as before $n - q = n - r - 1 \geq 1$.

(iv) We dispose of the case $m \geq r + 2$ and $n = r + 1$ through the following decomposition

$$S_A = \bigoplus_{i=1}^{r} \mathfrak{D}/_{\mathfrak{h}_i} \oplus \mathfrak{d} \oplus \mathfrak{D}^{n-r-1} = \bigoplus_{i=1}^{r} \mathfrak{D}/_{\mathfrak{h}_i} \oplus \mathfrak{d}$$

$$= \bigoplus_{i=1}^{r} \mathfrak{D}/_{\mathfrak{h}_i} \oplus \mathfrak{d} \oplus \mathfrak{D}^{1-1}$$

$$= S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{D})_{m \times r} \oplus S(\varnothing; \mathfrak{d})_{m \times 1}.$$

Clearly none of the summands vanishes and the first summand satisfies the empty part of condition (b) of the Seperated Divisor Theorem.

(v) If $m \geq r + 2$ and $n = r$, the module $S_A$ is a torsion module. As $m > h_{\mathfrak{D}}$ or $n > h_{\mathfrak{D}}$ we may invoke Lemma 4.28(ii), which proves the claim.

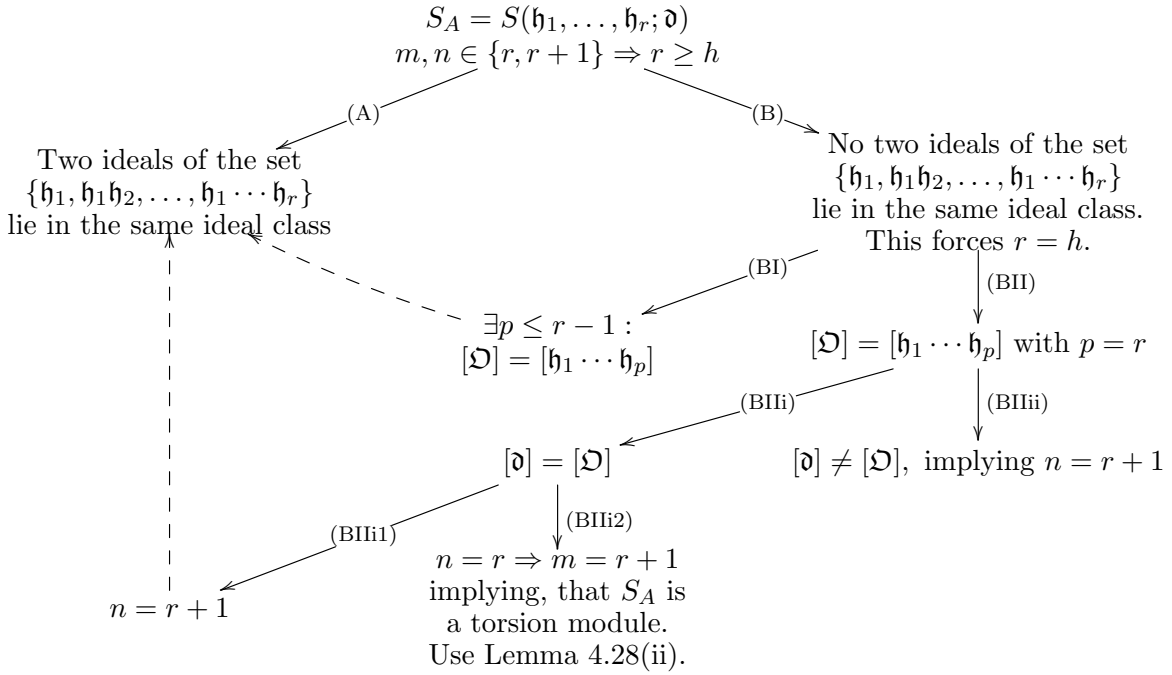$\square$

**Remark 4.30.** Note that the use of the assumption $m > h_{\mathfrak{O}}$ or $n > h_{\mathfrak{O}}$ in the previous two Lemmata 4.28(ii) and 4.29 were only necessary to handle torsion modules. Reviewing the proof of Lemma 4.28 one sees that it would have been sufficient to demand

$$m = 2 \wedge n = r = 1 \Rightarrow h_{\mathfrak{O}} = 1.$$

*Proof of Theorem 4.3.* Let $A$ be an $m \times n$ matrix of rank $r$. We show that, if $m$ or $n$ is greater $h := h_{\mathfrak{O}}$, then $A$ must be decomposable.
Cases $m \geq r + 2$ or $n \geq r + 2$ may be disposed of by the former Lemma 4.29.

We are thus left with $m, n \in \{r, r + 1\}$, which implies $r \geq h$. Let $S_A$ be given by $S_A = S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{d})_{m \times n}$. Since the proof requires a rather involved case-by-case analysis, we display the steps of the proof in a diagram. The arrow's label will indicate the specific case we are working on.

$$S_A = S(\mathfrak{h}_1, \ldots, \mathfrak{h}_r; \mathfrak{d})$$
$$m, n \in \{r, r + 1\} \Rightarrow r \geq h$$

(A)

(B)

Two ideals of the set
$\{\mathfrak{h}_1, \mathfrak{h}_1\mathfrak{h}_2, \ldots, \mathfrak{h}_1 \cdots \mathfrak{h}_r\}$
lie in the same ideal class

No two ideals of the set
$\{\mathfrak{h}_1, \mathfrak{h}_1\mathfrak{h}_2, \ldots, \mathfrak{h}_1 \cdots \mathfrak{h}_r\}$
lie in the same ideal class.
This forces $r = h$.

(BI)

(BII)

$\exists p \leq r - 1 :$
$[\mathfrak{D}] = [\mathfrak{h}_1 \cdots \mathfrak{h}_p]$

$[\mathfrak{D}] = [\mathfrak{h}_1 \cdots \mathfrak{h}_p]$ with $p = r$

(BIIi)

(BIIii)

$[\mathfrak{d}] = [\mathfrak{D}]$

$[\mathfrak{d}] \neq [\mathfrak{D}]$, implying $n = r + 1$

(BIIi1)

(BIIi2)

$n = r + 1$

$n = r \Rightarrow m = r + 1$
implying, that $S_A$ is
a torsion module.
Use Lemma 4.28(ii).

(A) Assume first that $\left[ \prod_{i=1}^{u} \mathfrak{h}_i \right] = \left[ \prod_{i=1}^{v} \mathfrak{h}_i \right]$ with some $u < v$. By multiplication with the inverses we find $[\mathfrak{D}] = \left[ \prod_{i=u+1}^{v} \mathfrak{h}_i \right]$. After renumbering of the $\mathfrak{h}_i$ we get

$$[\mathfrak{D}] = \left[ \prod_{i=1}^{p} \mathfrak{h}_i \right] \text{ with } p \leq r - 1. \tag{4.2}$$

Employing the Separated Divisor Theorem 4.24 we compute

$$
\begin{aligned}
S_A &\cong \bigoplus_{i=1}^{r} \mathfrak{D}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{D}^{n-r-1} \\
&= \bigoplus_{i=1}^{p} \mathfrak{D}/\mathfrak{h}_i \oplus \bigoplus_{i=p+1}^{r} \mathfrak{D}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{D}^{n-r-1} \\
&= \bigoplus_{i=1}^{p} \mathfrak{D}/\mathfrak{h}_i \oplus \bigoplus_{i=p+1}^{r} \mathfrak{D}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{D}^{n-p-(r-(p+1)+1)-1} \\
&= S(\mathfrak{h}_1,\ldots,\mathfrak{h}_p;\mathfrak{D})_{p\times p} \oplus S(\mathfrak{h}_{p+1},\ldots,\mathfrak{h}_r;\mathfrak{d})_{(m-p)\times(n-p)} \qquad (4.3) \\
&=: S_B \oplus S_C
\end{aligned}
$$

for suitable matrices $B, C$, the existence of which is guaranteed by the Separated Divisor Theorem. None of the summands vanishes, as $n - p \underset{(4.2)}{\geq} n - r + 1 \geq 1$, since $n \in \{r, r+1\}$; the same calculations hold true for $m$. Thus Lemma 4.28(i) evidences $A \sim \operatorname{diag}(B, C)$

The special cases of the Separated Divisor Theorem occur, if the number $r - p$ of ideals $\mathfrak{h}_i$ of the second summand in (4.3) equals $m - p$ or $n - p$, which implies $r \in \{m, n\}$. If $r = m$, then (a) of the Separated Divisor Theorem applied to $A$ shows

$$
[\mathfrak{d}] = \Big[\prod_{i=1}^{r} \mathfrak{h}_i\Big] = \Big[\prod_{i=1}^{p} \mathfrak{h}_i\Big]\Big[\prod_{i=p+1}^{r} \mathfrak{h}_i\Big]
$$

$$
\underset{(4.2)}{=} [\mathfrak{D}]\Big[\prod_{i=p+1}^{r} \mathfrak{h}_i\Big] = \Big[\prod_{i=p+1}^{r} \mathfrak{h}_i\Big],
$$

as demanded for the second term in (4.3).

When $r = n$, then case (b) of the Separated Divisor Theorem occurs and $S_A$ is a torsion module. We may set $[\mathfrak{d}] = [\mathfrak{D}]$ as a notational convention. The second term of (4.3) becomes

$$
S(\mathfrak{h}_{p+1},\ldots,\mathfrak{h}_r;\mathfrak{d})_{(m-p)\times(r-p)} = \bigoplus_{i=p+1}^{r} \mathfrak{D}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{D}^{-1} \cong \bigoplus_{i=p+1}^{r} \mathfrak{D}/\mathfrak{h}_i.
$$

Now $\bigoplus_{i=p+1}^{r} \mathfrak{D}/\mathfrak{h}_i$ is not degenerate as $p \leq r - 1$, thus the sum in (4.3) is a decomposition of $S_A$ into two modules.

(B) We are left with the case, where no two ideals

$$
\mathfrak{h}_1, \mathfrak{h}_1\mathfrak{h}_2, \ldots, \mathfrak{h}_1\cdots\mathfrak{h}_r \qquad (4.4)
$$

lie in the same ideal class. This together with $r \geq h$ implies $r = h$. Whence the ideals in (4.4) constitute a full set of representatives of the ideal class group.

(BI) The case where $[\mathfrak{h}_1 \mathfrak{h}_2 \cdots \mathfrak{h}_p] = [\mathfrak{O}]$ for some $p \leq r - 1$ can be treated analogously to (4.2) in case (A).

(BII) In the other case we have $[\mathfrak{h}_1 \mathfrak{h}_2 \cdots \mathfrak{h}_p] = [\mathfrak{O}]$, where $p = r$.

(BIIii) If $[\mathfrak{d}] \neq [\mathfrak{O}]$ and $[\mathfrak{h}_1 \cdots \mathfrak{h}_r] = [\mathfrak{O}]$. The ideal classes are determined by (4.4). As $S_A \cong \bigoplus_{i=1}^r \mathfrak{O}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{O}^{n-r-1}$ and as $\mathfrak{d}$ is not principal, we must have $n \neq r$ implying $n = r + 1$. Since (4.4) is a full system of representatives of ideal classes, there exists $q \in \mathbb{N}$, such that $[\mathfrak{d}] = [\mathfrak{h}_1 \cdots \mathfrak{h}_q]$. Note that $q \leq r - 1$, as $\mathfrak{h}_1 \cdots \mathfrak{h}_r$ is principal by assumption. We have $r < m$, since otherwise $[\mathfrak{d}] = \mathfrak{h}_1 \cdots \mathfrak{h}_r = [\mathfrak{O}]$ due to (a) of the Separated Divisor Theorem, which would contradict the assumption $[\mathfrak{d}] \neq [\mathfrak{O}]$. We review the inequalities we have gathered

$$q \leq r - 1 \text{ and } r < m \text{ implying } q < m, \qquad (4.5)$$
$$q + 1 \leq r < r + 1 = n.$$

We finish the proof by decomposing

$$
\begin{aligned}
S_A &\cong \bigoplus_{i=1}^r \mathfrak{O}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{O}^{n-r-1} \\
&= \bigoplus_{i=1}^q \mathfrak{O}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \bigoplus_{i=q+1}^r \mathfrak{O}/\mathfrak{h}_i \oplus \mathfrak{O}^{n-r-1} \\
&= \bigoplus_{i=1}^q \mathfrak{O}/\mathfrak{h}_i \oplus \mathfrak{d} \oplus \mathfrak{O}^{q+1-q-1} \oplus \bigoplus_{i=q+1}^r \mathfrak{O}/\mathfrak{h}_i \oplus \mathfrak{O} \oplus \mathfrak{O}^{n-q-1-(r-q)-1} \\
&= S(\mathfrak{h}_1, \ldots, \mathfrak{h}_q; \mathfrak{d})_{q \times (q+1)} \oplus S(\mathfrak{h}_{q+1}, \ldots, \mathfrak{h}_r; \mathfrak{O})_{(m-q) \times (n-q-1)} \\
&=: S_B \oplus S_C,
\end{aligned}
$$

for suitable matrices $B, C$. Due to the inequalities in (4.5) the matrix $C$ has a positive number of rows and columns. Now Lemma 4.28(i) implies $A \sim \mathrm{diag}(B, C)$.

(BIIi) Suppose that $\mathfrak{d}$ is principal, i.e. $[\mathfrak{d}] = [\mathfrak{O}]$. Considering the decomposition in (A) the second term of (4.3) degenerates to $S(\varnothing, \mathfrak{d})_{(m-p) \times (n-p)} = S(\varnothing, \mathfrak{d})_{(m-r) \times (n-r)}$. We need to treat the cases $n = r$ and $n = r + 1$ separately.

(BIIi1) If $n = r+1$, then $S(\varnothing, \mathfrak{d})_{(m-r) \times (n-r)} = \mathfrak{d}$. Hence the second module in (4.3) does not vanish and the decomposition of (A) is applicable.

(BIIi2) If $n = r$, then $m = r + 1$ follows, as we assumed $m > h = r$ or $n > h = r$. As $n = r$ the Separated Divisor Theorem yields that $S_A$ must be a torsion module. Invoking Lemma 4.28(ii) finishes the proof.

$\square$

*5*

<div style="background:#d3d3d3">

**Unit Equations**

</div>

This section prepares the tools necessary for dealing with questions concerning sums of units. Though we will solely be needing Theorem 5.5 in the progress, we take the time to survey the modern development of equations with units in a ring as solutions. Essentially all results are based on effective versions of Schmidt's subspace theorem [37] [38], a vast generalisation of the Thue-Siegel-Roth theorem in Diophantine approximation:

Denote by $\mathcal{O}$ the ring of all algebraic integers and let $\mathbf{x} = (x_1, \ldots, x_n) \in \mathbb{Q}^n$. Let $L_1, \ldots, L_n$ be linearly independent linear forms with coefficients from $\mathcal{O}$ mapping $\mathbb{Q}^n$ to $\mathbb{C}$. For every $\epsilon > 0$ there exist a finite number $V_1, \ldots, V_t$ of proper linear subspaces of $\mathbb{Q}^n$, such that the set of integer solutions to

$$|L_1(\mathbf{x}) \cdots L_n(\mathbf{x})| \leq ||\mathbf{x}||^{-\epsilon}, \quad \mathbf{0} \neq \mathbf{x} \in \mathbb{Z}^n \tag{5.1}$$

is contained in $V_1, \ldots, V_t$.

For the most recent development on quantitative results, i.e. effective upper bounds on the number of subspaces $t$, see the work of Evertse and Ferretti [39]. The latter authors also outline a result by Faltings and Wüstholz [40], who proved, that there exists a single, effectively computable, proper linear subspace of $\mathbb{Q}^n$, that contains almost all solutions to equation (5.1).

Turning to equations, where we require the solutions to be units, we take a top-down approach starting with a very general equation and successively specialising it. Let $\mathbf{x} = (x_1, \ldots, x_n) \in \mathbb{Z}^n$, $F \subseteq \mathbb{C}$ a field and $P_\ell \in F[X_1, \ldots, X_n]$. For $\mathbf{a}_\ell = (a_{\ell 1}, \ldots, a_{\ell n}) \in (F^*)^n$, we define $\mathbf{a}_\ell^{\mathbf{x}} = \prod_{i=1}^{n} a_{\ell i}^{x_i}$. The first type of equation we examine is a polynomial-exponential one:

$$\sum_{\ell=1}^{k} P_\ell(\mathbf{x}) \mathbf{a}_\ell^{\mathbf{x}} = 0. \tag{5.2}$$

A solution $\mathbf{x} \in \mathbb{Z}^n$ is called *degenerate* if

$$\exists I \subsetneqq \{1, \ldots, k\} : \sum_{\ell \in I} P_\ell(\mathbf{x}) \mathbf{a}_\ell^{\mathbf{x}} = 0.$$

Furthermore we say equation (5.2) fulfils the *Laurent-restriction*, if there are no non-trivial solutions to

$$\mathbf{a}_1^{\mathbf{z}} = \cdots = \mathbf{a}_k^{\mathbf{z}} \text{ in } \mathbf{z} \in \mathbb{Z}^n.$$

Much effort has been made to find effective bounds on the solutions of equation (5.2). Focusing on integer solutions, we introduce a theorem by Schlickewei and Schmidt.

**Theorem 5.1** ([41, 1.1$\mathcal{P}$]). *In the setting of (5.2), let $F$ be an algebraic number field of degree $d$ and. Define*

$$A = \sum_{\ell=1}^{k} \binom{n + \delta_\ell}{n} \text{ and } B = \max(n, A),$$

*where $\delta_\ell = \deg(P_\ell)$ is the total degree of $P_\ell$. If (5.2) satisfies the Laurent-restriction, then the number of non-degenerate solutions $\mathbf{x} \in \mathbb{Z}^n$ is bounded by*

$$N(d, B) := 2^{35B^3} d^{6B^2}.$$

Setting the length of the sum $k = 2$ in equation (5.2) Schmidt et al. [42] were able to eliminate the dependency on the number field's degree.

**Theorem 5.2** ([42]). *The equation*

$$\boldsymbol{a}^{\boldsymbol{x}} = P(\boldsymbol{x}) \text{ with } \boldsymbol{x} \in \mathbb{Z}^n,$$

*where $P \in \mathbb{C}[X_1, \ldots, X_n]$ and $\boldsymbol{a} = (a_1, \ldots, a_n) \in (\mathbb{C}^*)^n$ are multiplicatively independent has no more than*

$$\exp(B^{9B})$$

*solutions where $B = \binom{2 + \deg(P)}{2} + 1$, $\deg(P) > 0$.*

Next we specialise to a scalar version of (5.2). Consider the equation

$$\sum_{\ell=1}^{k} P_\ell(n) a_\ell^n = 0 \tag{5.3}$$

in the scalar unknown $n \in \mathbb{Z}$ with $P_\ell \in F[X], a_\ell \in F^*$, where $F \subseteq \mathbb{C}$. Note that equation (5.3) stems in fact from a linear recurrence sequence. In the next result due to Schmidt [43] the dependencies could again be reduced to the mere dependency on

the degree of the occurring polynomials. It can be viewed as an effective version of the Skolem-Mahler-Lech theorem [44] on linear recurrence sequences.

**Theorem 5.3** ([43]). *Let $F = \mathbb{C}$ and assume that $a_i a_j^{-1}$ is not a root of unity for $i \neq j$. Setting $t = \sum_{\ell=1}^{k} \big( \deg(P_\ell) + 1 \big)$ we have*

$$\#\{n \in \mathbb{Z} : \sum_{\ell=1}^{k} P_\ell(n) a_\ell^n = 0\} \leq \exp\big( \exp(t^{3t})\big).$$

The bound has recently been lowered to $e^{e^{70t}}$ by Amoroso [45, Theorem 1,2].

For the next step of specialisation we take the polynomials $P_\ell$ to be constant; Evertse, Schlickewei and Schmidt proved

**Theorem 5.4** ([8, Inequality (1.18)]). *In equation (5.3) let $F = \mathbb{C}$. Let $\boldsymbol{\alpha}, \boldsymbol{a} \in (\mathbb{C}^*)^k$ and assume that $(\mathbf{a.e}_i)(\mathbf{a.e}_j)^{-1}$ is not a root of unity for $i \neq j$, where $\mathbf{a.e}_i$ is the i-th component $\mathbf{a}$. We have*

$$|\{n \in \mathbb{Z} : \boldsymbol{\alpha}.\boldsymbol{a}^n = 0\}| \leq \exp(6k^{3k}).$$

The next theorem can be deduced from the latter, first stated as ineffective version by Evertse, Györy in [46]. The bound depends on the length of the sum and on a property of the group, we want the solutions to be contained in.

**Theorem 5.5** ([8]). *Let $\Gamma$ be a (discrete) subgroup of $(\mathbb{C}^*)^k$ of finite rank $r$. The equation*

$$\boldsymbol{\alpha}.\boldsymbol{a} = 1 \tag{5.4}$$

*in $\boldsymbol{a} \in \Gamma$ with coefficients $\boldsymbol{\alpha} \in (\mathbb{C}^*)^k$, has its number of non-degenerate solutions bound by*

$$M_k = \exp\big(6k^{3k}(r+1)\big).$$

For convenience we state a slight generalisation of the previous theorem as a system of equations.

**Corollary 5.6** ([46]). *Keep the setting of the former Theorem 5.5 and let $A$ be an $\ell \times k$ matrix over $\mathbb{C}^*$ and $0 \neq \boldsymbol{v} \in \mathbb{C}^\ell$. The number of non-degenerate solutions $\boldsymbol{x} \in \Gamma$ of*

$$A\boldsymbol{x} = \boldsymbol{v}$$

*is bounded by $\exp\big(6k^{3k}(r+1)\big)$.*

<div style="text-align: right">

# 6

</div>

<div style="text-align: right">

## $\omega$-good **rings**

</div>

## 6.1 Products and even rings

We start with some basic properties and definitions helpful for dealing with products of $\omega$-good rings. When dealing with sums of units, it is surprising that the notion of *even* rings as introduced in Chapter 1 does not seem to transcend Raphael's work [14], though it might be interesting to classify rings of algebraic integers with respect to the even-property. For convenience we extend the definition of even to number fields $K$, saying $K$ is $k$-even, if $\mathcal{O}_K$ satisfies the property.

**Proposition 6.1.** *The only $2$-even quadratic number fields are $\mathcal{O}_{-3}$ and $\mathcal{O}_5$. The representation of $1$ by sums of two units is given by*

$$1 = \frac{1}{2}\left(1 - \sqrt{-3}\right) + \frac{1}{2}\left(1 + \sqrt{-3}\right) \ \ and$$
$$1 = \frac{1}{2}\left(1 + \sqrt{5}\right) + \frac{1}{2}\left(1 - \sqrt{5}\right),$$
$$1 = \frac{1}{2}\left(-1 + \sqrt{5}\right) + \frac{1}{2}\left(3 - \sqrt{5}\right) = \frac{1}{2}\left(-1 - \sqrt{5}\right) + \frac{1}{2}\left(3 + \sqrt{5}\right)$$

*Proof.* We claim, that $\mathcal{O}_d$ is not 2-even, if $d \not\equiv 1 \pmod 4$:
The generic element is of the form $a + b\sqrt{d}$ with $a, b \in \mathbb{Z}$. We require

$$1 = (a + b\sqrt{d}) + (e + f\sqrt{d}),$$

where both factors are units in $\mathcal{O}_d$. A small conversion yields $e + f\sqrt{d} = (1-a) - b\sqrt{d} \in \mathcal{O}_d^*$. Considering the norm we conclude from

$$\begin{aligned}
\pm 1 &= a^2 - b^2 d^2 \\
&= (1-a)^2 - b^2 d^2
\end{aligned} \tag{6.1}$$

that $2a - 1 = 0$, which is impossible.

Turning to $d \equiv 1 \pmod 4$, the generic element takes the form $a + b\frac{1+\sqrt{d}}{2}$ with $a, b \in \mathbb{Z}$. Using the same ansatz (6.1) as before we are led to the equations

$$
\begin{aligned}
\pm 1 &= \mathcal{N}\left(a + b\frac{1+\sqrt{d}}{2}\right) & &= a^2 + ab + \frac{1-d}{4}b^2 \\
&= \mathcal{N}\left((1-a) - b\frac{1+\sqrt{d}}{2}\right) & &= (1-a)^2 - (1-a)b + \frac{1-d}{4}b^2.
\end{aligned}
$$

We have $0 = 2a + b - 1$, which shows $2a + b = 1$. As

$$
\begin{aligned}
1 &= \left|\mathcal{N}(a + b\frac{1+\sqrt{d}}{2})\mathcal{N}(1 - a - b\frac{1+\sqrt{d}}{2})\right| \\
&= \left|\frac{1}{16}\left((2a + b - 2)^2 - b^2 d\right)\left((2a + b)^2 - b^2 d\right)\right|,
\end{aligned}
$$

we see that $1 = \frac{1}{16}\left|(1 - b^2 d)^2\right|$. This implies $1 - b^2 d = \pm 4$, thus $b = \pm 1$, $d \in \{-3, 5\}$.

The problem of finding $a$, respectively of showing that the solutions are indeed as claimed, is easily settled.       □

**Remark 6.2.** Let $K$ be a number field of even degree $k$. Let $\epsilon \in \mathcal{O}_K$ have degree $k$ and suppose

$$|\operatorname{Tr}_K(\epsilon)| = |\mathcal{N}_K(\epsilon)| = 1.$$

Then $\mathcal{O}_K$ is $k$-even.

*Proof.* Note first that an $\epsilon$ of degree smaller $k$ can not have $|\operatorname{Tr}_K(\epsilon)| = 1$. From the condition regarding the trace we obtain

$$\pm 1 = \operatorname{Tr}_K(\epsilon) = \sum_{i=1}^{k} \epsilon_i,$$

where the $\epsilon_i$'s are the conjugates of $\epsilon$ in $\mathcal{O}_K$. The condition on the norm of $\epsilon$ implies that all its conjugates are units as well, whence follows the claim.       □

**Remark 6.3.** Anticipating Proposition 6.9 of the next section, we have that $\mathcal{O}_5$ and $\mathcal{O}_{-3}$ are not only 2-even but also $\omega$-good. Using the previous computations of Proposition 6.1 we construct new examples of 2-even rings.

**Example 6.4.**

(i) Employing Lemma 1.6(i) it is clear, that the extension of a $k$-even ring is again $k$-even. Using this fact it is easy to construct extensions of $\mathbb{Q}(\sqrt{5})$ being 2-even: As the Legendre symbol $\left(\frac{5}{3}\right)$ is equal to $-1$, we see that 3 is inert in $\mathcal{O}_5$, hence prime. This shows that the polynomial $x^d - 3 \in \mathcal{O}_5[x]$ is Eisensteinian for all $d \geq 2$. Therefore a solution $c_d$ to the equation $x^d - 3 = 0$ has degree $d$ over

$\mathbb{Q}(\sqrt{5})$. It follows that $\mathbb{Q}(\sqrt{-5}, c_d)$ has degree $2d$ over $\mathbb{Q}$ and is 2-even, where $d \geq 2$ may be chosen arbitrarily.

(ii) Next we derive a result concerning the property even for cyclotomic field. Let $n \in \mathbb{N}$ not be a power of 2, denote by $p$ an arbitrary prime factor of $n$ greater 2. Then the $n$-th cyclotomic field is $(p-1)$-even: let $\zeta_p$ denote a primitive $p$-th root of unity, then

$$0 = \sum_{i=0}^{p-1} \zeta^i \text{ and therefore } 1 = -\sum_{i=1}^{p-1} \zeta^i.$$

Thus $\mathbb{Q}(\zeta_p)$ is $(p-1)$-even, and as $\mathcal{O}_{\zeta_p}$ injects into $\mathcal{O}_{\zeta_n}$, we see that the number field $\mathbb{Q}_{\zeta_n}$ of degree[1] $\varphi(n)$ over $\mathbb{Q}$ is also $(p-1)$-even.

**Definition 6.5.** Let $\mathcal{O}$ be a ring of algebraic integers with quotient field $K$. Let the $k$-good elements of $\mathcal{O}$ be denoted by $\Sigma_k$. We call $\mathcal{O}$ *active*, if $\mathcal{O} \neq \mathbb{Z}$ and the inertia degree of all prime ideals dividing $(2) = 2\mathcal{O}$ is greater one. In particular this is fulfilled, if 2 is prime in $\mathcal{O}$. On the other hand, the condition is violated, if $(2)$ splits in $\mathcal{O}$ or $(2)$ is totally ramified in $\mathcal{O}$.

**Theorem 6.6.** *Let $\mathcal{O}$ be an active ring of algebraic integers. If $\mathcal{O}$ is $\omega$-good, then $\mathcal{O}$ is even.*

*Proof.* Suppose $\mathcal{O}$ is $\omega$-good, but $\mathcal{O}$ is not even, i.e. odd. Define

$$u : \mathcal{O} \to \mathbb{F}_2$$
$$\alpha \mapsto \begin{cases} 0 & \alpha \in \bigcup_{i \geq 1} \Sigma_{2i} \\ 1 & \alpha \in \bigcup_{i \geq 1} \Sigma_{2i-1} \end{cases},$$

The map is well-defined: Clearly $\mathcal{O}$ is the union of $\bigcup_{i \geq 1} \Sigma_{2i}$ and $\bigcup_{i \geq 1} \Sigma_{2i-1}$ as $\mathcal{O}$ is $\omega$-good. To see that the union is in fact disjoint, suppose there is $\alpha \in \bigcup_{i \geq 1} \Sigma_{2i} \cap \bigcup_{i \geq 1} \Sigma_{2i-1}$. Then there exist $k \in 2\mathbb{N}, \ell \in \mathbb{N} \backslash 2\mathbb{N}$ and units $\epsilon_i, \eta_i \in \mathcal{O}$, such that $\alpha = \sum_{i=1}^{k} \epsilon_i = \sum_{i=1}^{\ell} \eta_i$. Now $\eta_\ell = \sum_{i=1}^{k} \epsilon_i - \sum_{i=1}^{\ell-1} \eta_i$ is a sum of $k + \ell - 1 \in 2\mathbb{N}$ units contradicting the assumption that $\mathcal{O}$ is odd.

It is easy to check that $u$ constitutes a ring homomorphism. Set $\mathfrak{w} := \ker(u)$, then $\mathcal{O}/\mathfrak{w} \cong \mathbb{F}_2$, which shows that $\mathfrak{w}$ is a prime ideal of $\mathcal{O}$ lying over the rational prime 2 featuring inertia degree 1 - a contradiction.

$\square$

---

[1] $\varphi$ denotes Euler's totient function.

**Corollary 6.7.** *Let $d > 0$ be a positive, squarefree integer contained in the sequence*

$$a_n = \frac{5}{2}\left((-1)^n + 1\right) + n\left(n - (-1)^n + 3\right).$$

*Then $\mathcal{O}_d$ is even.*

*Proof.* A short calculation evidences that the positive integers $a_n$ are precisely those positive integers $a$ fulfilling both of the following properties:

- $a + 4$ or $a - 4$ is a perfect square

- $a \equiv 5 \pmod 8$.

Anticipating Proposition 6.9, we see that $d$ satisfies condition (ii). Hence $O_d$ is $\omega$-good. Moreover it is a basic fact in algebraic number theory that 2 is inert in $\mathcal{O}_d$, if and only if $d \equiv 5 \pmod 8$. Thus $O_d$ meets the requirements of Theorem 6.6, proving that $\mathcal{O}_d$ is even. $\qquad\square$

 

The following corollary gives a sufficient condition for an algebraic ring of integers $\mathcal{O}$ to satisfy $u(\mathcal{O}) = \infty$, which also extends to subrings.

**Corollary 6.8.** *Let $\mathcal{O}$ be odd and 2 prime in $\mathcal{O}$. Then for any subfield $M \subseteq K \neq \mathbb{Q}$, we have $u(\mathcal{O}_M) = \infty$.*

*Proof.* It is easy to see that the properties *odd* and *2 being inert* are inherited by any subfield. By the negation of Theorem 6.6 we obtain, that any odd algebraic ring of integers $\mathcal{O}$, satisfying that 2 is inert in $\mathcal{O}$, in particular implying that $\mathcal{O}$ is active, has unit sum number $\infty$. Hence follows the claim.

$$\square$$

## 6.2  The unit sum number of number fields of small degree

Before turning to a general result, we observe some completely solved instances of low-degree number fields. The quadratic case has been investigated by Belcher as early as 1974.

**Proposition 6.9** ([3, Theorem 7], cf. [47, Lemma 1])**.** *Let $K = \mathbb{Q}(\sqrt{d})$, $d \in \mathbb{Z}$ squarefree. Then the ring of algebraic integers $\mathcal{O}_d$ is $\omega$-good, if and only if one of the following conditions hold:*

  *(i) $d > 0, d \not\equiv 1 \pmod 4$ and $d + 1$ or $d - 1$ is a perfect square*

  *(ii) $d > 0, d \equiv 1 \pmod 4$ and $d + 4$ or $d - 4$ is a perfect square*

  *(iii) $d \in \{-1, -3\}$*

*Proof.* We first attend to the case, where $\mathbb{Q}(\sqrt{d})$ is an imaginary quadratic field, i.e. $d < 0$. The unit structure of imaginary quadratic number fields is fully determined[2]:

$$\mathcal{O}_d{}^* = \begin{cases} \{\pm 1, \pm i\} & d = -1 \\ \{\pm 1, \pm \zeta, \pm \zeta^2\} & d = -3 \\ \{\pm 1\} & \text{else} \end{cases},$$

where $\zeta$ denotes a third primitive root of unity. It is readily observed that $\mathbb{Q}(\sqrt{-1})$ and $\mathbb{Q}(\sqrt{-3})$ are the only $\omega$-good imaginary quadratic extensions of the rationals.

We move on to consider the case $d > 0$.

"$\Leftarrow$"

(i) Let $d \not\equiv 1 \pmod 4$ and $d = a^2 \pm 1$ for some $a \in \mathbb{Z}$. In this case $\{1, \sqrt{d}\}$ constitutes an integral basis for $\mathbb{Q}(\sqrt{d})$. We have

$$(-a + \sqrt{d})(a + \sqrt{d}) = d - a^2 = \pm 1,$$

showing that $-a + \sqrt{d} \in \mathcal{O}_d^*$. Thus $\sqrt{d} = a + (-a + \sqrt{d})$ is a sum of units, as $a = \underbrace{1 + \cdots + 1}_{a\text{-times}}$. Therefore $u(\mathbb{Q}(\sqrt{d})) = \omega$.

(ii) If $d \equiv 1 \pmod 4$ and $d = a^2 \pm 4$ for some $a \in \mathbb{Z}$, then $a$ is odd and $\{1, \frac{1+\sqrt{d}}{2}\}$ forms an integral basis. We calculate

$$\left(\frac{a-1}{2} + \frac{1+\sqrt{d}}{2}\right)\left(\frac{a-1}{2} + 1 - \frac{1+\sqrt{d}}{2}\right) = \frac{1}{4}(a^2 - d) = \pm 1,$$

observing that both factors are contained in $\mathcal{O}_d$ and are therefore units. As

$$\frac{1+\sqrt{d}}{2} = \frac{1-a}{2} + \left(\frac{a-1}{2} + \frac{1+\sqrt{d}}{2}\right),$$

the first factor being an integer in $\mathbb{Z}$, we see that $\mathbb{Q}(\sqrt{d})$ is $\omega$-good.

"$\Rightarrow$"

We have yet to show that for $d \in \mathbb{N}$ squarefree, not as in (i) or (ii) the unit sum number of $\mathbb{Q}(\sqrt{d})$ is infinite. Given an integral basis $\{1, \delta\}$, define the additive map

$$\pi : \mathbb{Z}[\delta] \mapsto \mathbb{Z}$$
$$\pi(r + s\delta) = s,$$

where $r, s \in \mathbb{Z}$. Dirichlet's unit theorem implies that all units of a real quadratic number field can be written as $\pm \eta^k$, where $\eta$ denotes a fundamental unit of $\mathbb{Q}(\sqrt{d})$.

---

[2]For a proof see [35, Proposition 3.1.8].

Without loss of generality we may set $\eta = a + p\delta$, with $p > 0$, as $-\eta$ is a fundamental unit as well. We find

$$\pi(\pm\eta^k) = \pi(\pm(a + p\delta)^k),$$

which is divisible by $p$ due to the Binomial theorem. Thus $\pi(\mathcal{O}^\omega) \subseteq p\mathbb{Z}$, where $\mathcal{O}^\omega$, as introduced in Chapter 1, denotes the subring of $\mathcal{O}_d$ containing all $\omega$-good elements. Suppose now $\mathcal{O}_d$ is $\omega$-good, then $\delta \in \mathcal{O}^\omega$ and $\pi(\delta) = 1 \in p\mathbb{Z}$, hence $p = 1$. Therefore $\pi(\eta) = 1$ and $\eta = a + \delta$. Clearly $\eta$ has norm $\mathcal{N}(\eta) = \pm 1$. First let $d \not\equiv 1 \pmod 4$, which entails $\delta = \sqrt{d}$, then

$$\mathcal{N}(a + \delta) = a^2 - d = \pm 1 \text{ as in case (i).}$$

Suppose $d \equiv 1 \pmod 4$ entailing $\delta = \frac{1+\sqrt{d}}{2}$, then

$$\mathcal{N}(a + \delta) = a^2 + a + \frac{1 - d}{4} = \pm 1.$$

As $4a^2 + 4a + 1 = (2a + 1)^2$ is a perfect square and we are led to case (ii).    $\square$

Pure cubic number fields, i.e. fields of the form $\mathbb{Q}(\sqrt[3]{d})$ with $d \in \mathbb{Z}$ cubefree, have been examined by Tichy and Ziegler:

**Proposition 6.10** ([48, Theorem 2]). *Let $K = \mathbb{Q}(\sqrt[3]{d})$ be a pure cubic field, $\mathcal{O}_K$ is $\omega$-good, if and only if one of the two cases hold*

   *(i) $d$ is squarefree, $d \not\equiv \pm 1 \pmod 9$ and $d + 1$ or $d - 1$ is a perfect cube*

   *(ii) $d = 28$.*

Further results have been obtained by Filipin, Tichy and Ziegler [49, Theorem 1.1] for pure quartic complex fields $\mathbb{Q}(\sqrt[4]{d})$, with $\mathbb{Z} \ni d < 0$ and $d \neq -4$. We remark that all of these three special cases have in common, that due to Dedekind's unit theorem their unit group-rank is equal to one, i.e. there is only one fundamental unit.

There is yet another type of number field, which we easily recognize as being $\omega$-good, namely cyclotomic fields. The next theorem warrants an integral basis consisting solely of units.

To facilitate the readability of the proof, we use an auxiliary lemma to settle some calculations.

**Lemma 6.11.** *Let $p > 2$ be prime and $\alpha \in \mathbb{N}$, then*

$$\sum_{j=0}^{p-1}(1 + x^\alpha)^j \equiv x^{\alpha(p-1)} \pmod p.$$

*Proof.* Using the Binomial theorem and exchanging the order of summation, we obtain

$$\sum_{j=0}^{p-1} (1 + x^\alpha)^j = \sum_{j=0}^{p-1} \sum_{k=0}^{j} \binom{j}{k} x^{\alpha k} = \sum_{j=0}^{p-1} \left( \sum_{k=0}^{p-1} \binom{k}{j} \right) x^{\alpha j}.$$

Set

$$\chi_n(j) = \sum_{k=0}^{n-1} \binom{k}{j} \quad \text{and} \quad \gamma_n(j) = \frac{n-j}{j+1} \binom{n}{j}$$

The proof is complete, if we are able to deduce $\chi_n(j) = \gamma_n(j)$ for arbitrary $n, j \in \mathbb{N}_0$, as this leads to

$$\chi_p(j) = \gamma_p(j) \equiv \delta_{j,p-1} \pmod{p} \text{ for } 0 \le j \le p-1.$$

It is easy to observe that $\chi_1 = \gamma_1$. Suppose for a fixed $n$ that $\chi_n = \gamma_n$, then

$$\chi_{n+1}(j) = \sum_{k=j}^{n} \binom{k}{j} = \gamma_{n-1}(j) + \binom{n}{j} = \frac{n-j}{j+1} \binom{n}{j} + \binom{n}{j}.$$

Multiplying the following equation by $\frac{1}{1+j}$ shows that $\chi_{n+1} = \gamma_{n+1}$:

$$(n-j) \binom{n}{j} + (1+j) \binom{n}{j} = (n+1) \binom{n}{j} = \frac{(n+1)n!}{j!(n-j)!} =$$
$$= \frac{(n-j+1)(n+1)!}{j!(n-j+1)!} = (n-j+1) \binom{n+1}{j}.$$

$\square$

**Theorem 6.12** ([35, Theorem 4.27, Theorem 2.20]). *Let $m \in \mathbb{N}$, $\zeta_m$ an $m$-th primitive root of unity. Then the $m$-th cyclotomic field $K_m = \mathbb{Q}(\zeta_m)$ has an integral basis $\{1, \zeta_m, \zeta_m^2, \ldots, \zeta_m^{\varphi(m)-1}\}$, where $\varphi$ denotes Euler's totient function. In particular $\mathcal{O}_{\mathbb{Q}(\zeta_m)} = \mathbb{Z}[\zeta_m]$.*[3]

*Proof.* We will proof the theorem only for $m = p^n$ with some $n \in \mathbb{N}$, $p$ prime. Set $\Phi(x) = \frac{x^{p^n}-1}{x^{p^{n-1}}-1}$. By virtue of L'Hospital's rule we evaluate $\Phi(x)$ at some points for later use:

$$\Phi(1) = \lim_{x \to 1} \frac{x^{p^n}-1}{x^{p^{n-1}}-1} = \lim_{x \to 1} \frac{p^n x^{p^n-1}}{p^{n-1} x^{p^{n-1}-1}} = p,$$

furthermore $\Phi(\zeta_m) = 0$. Consider the shifted polynomial $\tilde{\Phi}(x) = \Phi(x+1)$.

---

[3] A ring of algebraic integers generated by powers of a single element is called *monogenic*. For a short introduction see [35, §2.6], for a proof that there are monogenic number fields of every signature see [50].

For $j \in \mathbb{N}$ we have

$$(1+x)^{jp^{n-1}} \equiv (1+x^{p^{n-1}})^j \pmod{p},$$

which is easily evidenced via induction. Thus using geometric series and applying the former lemma, we get

$$\tilde{\Phi}(x) = \frac{(x+1)^{p^n} - 1}{(x+1)^{p^{n-1}} - 1} \equiv \frac{(1+x^{p^{n-1}})^p - 1}{(1+x^{p^{n-1}}) - 1}$$

$$= \sum_{j=0}^{p-1} (1+x^{p^{n-1}})^j \equiv x^{p^{n-1}(p-1)} \pmod{p}.$$

This computation together with $\tilde{\Phi}(0) = \Phi(1) = p$ reveals that $\tilde{\Phi}$ is $p$-Eisensteinian, hence $\Phi$ is irreducible and thus the minimal polynomial of $\zeta_m$. Regarding the degree we obtain $[K : Q] = p^{n-1}(p-1) = \varphi(p^n) = \varphi(m)$. Using the well-known formula

$$d_K(\zeta_m) = \pm \mathcal{N}(\Phi'(\zeta_m)), \tag{6.2}$$

and noting that $\tilde{\Phi}$ is the minimal polynomial of $\zeta_m - 1$, we arrive at $d_K(\zeta_m) = d_K(\zeta_m - 1)$. [4]

As the conjugates of $\zeta_m$ form a $\mathbb{Q}$-basis of $K$, we must have $d_K | d_K(\zeta_m) = d_K(\zeta_m - 1)$; the quotient is called the index of $\zeta_m$ in $K$. The proof is finished, if we verify that $d_K(\zeta_m)$ equals a power of $p$ and that the index is not divisible by $p$, whence $d_K = d_K(\zeta_m)$. After some calculations one arrives at

$$\Phi'(\zeta_m) = \frac{p^n}{\zeta_m(\zeta_m^{p^{n-1}} - 1)}.$$

Noting that $\mathcal{N}(\zeta_m) = \Phi(0) = 1$ and $\zeta_m^{p^{n-1}} = \zeta_p$ is some $p$-th primitive root of unity, formula (6.2) warrants

$$\pm d_K(\zeta_m) = \frac{p^{n\varphi(m)}}{\mathcal{N}(\zeta_p - 1)}.$$

By transitivity of the norm

$$\mathcal{N}_{K/\mathbb{Q}}(\zeta_p - 1) = \mathcal{N}_{\mathbb{Q}(\zeta_p)/\mathbb{Q}}\big(\mathcal{N}_{K/\mathbb{Q}(\zeta_p)}(\zeta_p - 1)\big)$$

$$= \mathcal{N}_{\mathbb{Q}(\zeta_p)/\mathbb{Q}}\big((\zeta_p - 1)^{\varphi(m)/(p-1)}\big)$$

$$= \mathcal{N}_{\mathbb{Q}(\zeta_p)/\mathbb{Q}}(\zeta_p - 1)^{p^{n-1}}.$$

Now

$$\mathcal{N}_{\mathbb{Q}(\zeta_p)/\mathbb{Q}}(\zeta_p - 1) = \prod_{j=1}^{p-1}(\zeta_p^j - 1) = (-1)^{p-1} \lim_{x \to 1} \frac{x^p - 1}{x - 1} = (-1)^{p-1} p,$$

---

[4]Correction in Narkiewicz's proof: $..\zeta_q - 1$ is the root of $F..$

hence $d_K(\zeta_m)$ equals some power of $p$ as claimed.

Since the minimal polynomial $\tilde{\Phi}(x) = x^{\varphi(m)} + a_{\varphi(m)-1}x^{\varphi(m)-1} + \cdots + a_0$ of $\zeta_m - 1$ is $p$-Eisensteinian, we have

$$\frac{1}{p}(\zeta_m - 1)^{\varphi(m)} \in \mathbb{Z} \text{ and } p^2 \nmid \mathcal{N}_{K/\mathbb{Q}}(\zeta_m - 1).$$

Suppose now the index of $d_K(\zeta_m - 1)$ is divisible by $p$. There exists $\mu \in \mathcal{O}_K$, $b_i \in \mathbb{Z}$, such that $p\mu = \sum_{i=0}^{\varphi(m)-1} b_i(\zeta_m - 1)^i$, where not all $b_i \in \mathbb{Z}$ are divisible by $p$. Letting $j$ be the minimal index with $p \nmid b_j$, we see

$$\eta := \frac{1}{p}\sum_{i=j}^{\varphi(m)-1} b_i(\zeta_m - 1)^i = \mu - \frac{1}{p}\sum_{i=0}^{j-1} b_i(\zeta_m - 1)^i \in \mathcal{O}_K.$$

Whence also $\beta := \frac{b_j}{p}(\zeta_m - 1)^{\varphi(m)-1}$, which is equal to

$$\eta(\zeta_m - 1)^{\varphi(m)-j-1} - \frac{(\zeta_m - 1)^{\varphi(m)}}{p}\big(b_{j+1} + b_{j+2}(\zeta_m - 1) + \cdots + b_{\varphi(m)-1}(\zeta_m - 1)^{\varphi(m)-j-2}\big)$$

is an algebraic integer. This leads to

$$p^{\varphi(m)}\mathcal{N}_{K/\mathbb{Q}}(\beta) = \mathcal{N}_{K/\mathbb{Q}}(p\beta) = \mathcal{N}_{K/\mathbb{Q}}\big(b_j(\zeta_m - 1)^{\varphi(m)-1}\big) = b_j^{\varphi(m)}\mathcal{N}_{K/\mathbb{Q}}(\zeta_m - 1)^{\varphi(m)-1},$$

and $p|b_j$ - a contradiction. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\square$

**Remark 6.13.** As we have seen, the integral bases of cyclotomic fields consist of units. This immediately implies that all cyclotomic fields are $\omega$-good. Clearly much more effort is needed to prove that a field is $\omega$-good in comparison to proving that it satisfies the even-property - confer Example 6.4(ii).

## 6.3 The unit sum number of number fields

The most important result regarding the unit sum number problem for algebraic number fields is due to Jarden and Narkiewciz. However, no general criteria are known, whether $\mathcal{O}_K$ is $\omega$-good or not.

**Definition 6.14.** An arithmetic progression of length $r \in \mathbb{N}$ in a commutative ring $R$ is a sequence $\{a_i\}_{i=1}^r \subseteq R$, satisfying $a_i = a_1 + id$ for all $i \in \{2, \ldots, r\}$ and a fixed $d \in R$.

**Theorem 6.15** ([7, Theorem 1]). *In a ring of algebraic integers $\mathcal{O}$ every arithmetic progression of k-good elements, $k \in \mathbb{N}$, has finite length. In particular the unit sum number of an algebraic number field $K$ is either $\omega$ or $\infty$.*

**Remark 6.16.** The theorem leads to an immediate improvement of Proposition 6.9, as it shows that the identified quadratic, $\omega$-good extensions $\mathcal{O}_K$ are not only $\omega$-good, but fulfil $u(\mathcal{O}_K) = \omega$. Similarly considering Remark 6.13 we deduce that cyclotomic fields have unit sum number $\omega$.

To proof the theorem we employ Theorem 5.5 and van der Waerden's theorem on *arithmetic progressions* in $\mathbb{Z}$.

### 6.3.1 Van der Waerden's theorem

**Theorem 6.17** ([9], cf. [51]). *We have $\forall k, \ell \in \mathbb{N} : \exists n = n(k, \ell) \in \mathbb{N}$, such that a partition of $\{1, \ldots, n\}$ into $k$ subsets, guarantees the existence of an arithmetic progression of length $\ell$ within one of these $k$ subsets.*

The smallest possible number $n(k, l)$, for which the theorem holds, is called *van der Waerden number*. Though still five-fold exponential the best general bound on $n(k, \ell)$ is due to Gowers [52, Theorem 18.2], who deduced

$$n(k, \ell) \leq 2^{2^{k^{2^{2^{(\ell+9)}}}}}.$$

A much easier accessible result was given by Berlekamp [53] about 30 years earlier, who proved for primes $t$, that $n(k, t) \geq t2^t$. There also exist exact result for small $k$. We provide a table of all exact values known to date:

| vdW. number | Value | Author |
|---|---|---|
| n(2,3) | 9 | Chvátal [54] |
| n(2,4) | 35 | Chvátal [54] |
| n(2,5) | 178 | Stevens and Shantaram [55] |
| n(2,6) | 1132 | Kouril and Paul [56] |
| n(3,3) | 27 | Chvátal [54] |
| n(3,4) | 293 | Kouril [57] |
| n(4,3) | 76 | Beeler and O'Neil [58] |

To prove the theorem we rely on an auxiliary lemma, which employs multi-arithmetic progressions. To simplify notation for $j < k \in \mathbb{N}$ we write $[j, k]$ to denote the set $\{j, j+1, \ldots, k\}$.

**Lemma 6.18** ([9, Lemma 2.1]). *Suppose van der Waerden's theorem holds for a certain $\ell \geq 2$ and $\forall k \in \mathbb{N}$. Then for all $k, m \in \mathbb{N}$ there is $N(k, m, \ell) \in \mathbb{N}$ with the following property: Set $\Delta = [1, N(k, m, \ell)]$ and let $\rho$ denote an arbitrary surjective function $\rho : \Delta \to [1, k]$. Then there exists a function*

$$f : [0, \ell]^m \to \Delta$$

$$f(i_1, \ldots, i_m) = a + \sum_{\nu=1}^{m} i_\nu d_\nu$$

*with suitable $a \in \mathbb{N}, d_\nu \in \mathbb{N}$, such that*

$$(\rho \circ f)(i_1, \ldots, i_s, j_{s+1}, \ldots, j_m) = (\rho \circ f)(0, \ldots, 0, j_{s+1}, \ldots, j_m)$$

*for all $j_{s+1}, \ldots, j_m$, whenever $i_1, \ldots, i_s \in [0, \ell-1]$. For convenience call such functions $m$-long.*

*Proof.* We induct on $m$, in the case of $m = 1$ setting

$$N(k, 1, \ell) = 2n(k, \ell),$$

where $n(k, \ell)$ is the van der Waerden number for $k$ and $\ell$ as before.
Let

$$\Delta_1 = [1, n(k, \ell)] \text{ and } \Delta_2 = [n(k, \ell) + 1, 2n(k, \ell)],$$

hence $\Delta = [1, N(k, 1, \ell)]$ is their union. Let $I \subseteq [1, k]$ denote the image of $\rho_{|\Delta_1}$. As $|I| \leq k$ we invoke van der Waerden's theorem to find $t \in I$ such that there exists a progression $a + id, 0 \leq i < \ell$, in $\rho^{-1}(t) \subseteq \Delta_1$. As $a + (\ell-1)d \leq n(k, l)$ and $d \leq n(k, \ell)$, we have $a + \ell d \in \Delta$. Therefore setting $f$ to be the progression $a + id, 0 \leq i \leq \ell$ in $\Delta$, we find $f([0, \ell-1]) \subseteq \rho^{-1}(t)$. Thus $\forall i \in [0, \ell-1] : (\rho \circ f)(i) = (\rho \circ f)(0)$ as required.

Assume now the lemma holds for a certain $m$. We define

$$\begin{aligned}
q :=& N(k, m, \ell)\\
N(k, m+1, \ell) :=& 2n(k^q, \ell) + N(k, m, \ell)\\
\Delta_1 :=& [1, 2n(k^q, \ell)]\\
\Delta_2 :=& [2n(k^q, \ell) + 1, N(k, m+1, \ell)],
\end{aligned}$$

and hence $\Delta = \Delta_1 \cup \Delta_2 = [1, N(k, m+1, \ell)]$. Define an equivalence relation $\sim$ on $\Delta_1$ by

$$\begin{aligned}
x \sim y \Leftrightarrow \rho(x) =& \rho(y) \wedge\\
\rho(x+1) =& \rho(y+1) \wedge\\
&\vdots\\
\rho(x+q-1) =& \rho(y+q-1).
\end{aligned}$$

The number of possible equivalence classes of $\sim$ is $k^q$, since $\rho$ maps to $[1, k]$ and thus each of the $q$ conditions evaluates to an integer in $[1, k]$. As by definition $|\Delta_1| = 2n(k^q, \ell)$ van der Waerden's theorem grants the existence of an arithmetic progression $a + id, 0 \le i < \ell$, which is contained in one of the $k^q$ equivalence classes of $\sim$ in $\Delta_1$. Equivalently we may write

$$a \sim a + d \sim \cdots \sim a + (\ell - 1)d.$$

This and the definition of $\sim$ warrants the following set of conditions

$$
\begin{aligned}
\rho(a) &= & \rho(a + d) &= \ldots & &= \rho(a + (\ell - 1)d) \\
\rho(a + 1) &= & \rho(a + d + 1) &= \ldots & &= \rho(a + (\ell - 1)d + 1) \\
&\vdots & &\vdots & &\vdots \\
\rho(a + q - 1) &= & \rho(a + d + q - 1) &= \ldots & &= \rho(a + (\ell - 1)d + q - 1),
\end{aligned}
$$

which may be stated as

$$\forall c \in [a, a + q) : \forall i \in [0, \ell) \text{ it is true that } \rho(c + id) = \rho(c). \tag{6.3}$$

As $[a, a + q) \subseteq \Delta$ has length $q = N(k, m, \ell)$ the induction hypothesis yields an $m$-long function $g : [0, \ell]^m \to [a, a + q)$. Define

$$
\begin{aligned}
f &: [0, \ell]^{m+1} \to \Delta \\
f(i_0, \ldots, i_m) &= i_0 d + g(i_1, \ldots, i_m).
\end{aligned}
$$

For $i_0, \ldots, i_s < \ell$ and $s \in \mathbb{N}$ we have

$$
\begin{aligned}
(\rho \circ f)(i_0, i_1, \ldots, i_s, j_{s+1}, \ldots, j_m) &= \rho\big(i_0 d + g(i_1, \ldots, i_s, j_{s+1}, \ldots, j_m)\big) \\
&= (\rho \circ g)(i_1, \ldots, i_s, j_{s+1}, \ldots, j_m) \\
&= (\rho \circ g)(0, \ldots, 0, j_{s+1}, \ldots, j_m) \\
&= (\rho \circ f)(0, 0, \ldots, 0, j_{s+1}, \ldots, j_m)
\end{aligned}
$$

for all $j_{s+1}, \ldots, j_m$, where the conversion from the first line to the second line is valid due to (6.3). $\qquad\square$

*Proof of 6.17.* The theorem holds trivially for arbitrary $k$ and $\ell = 2$, in which case we could set $n = k + 1$. Suppose the theorem holds for a certain $\ell \in \mathbb{N}$. This enables us to invoke the previous Lemma 6.18. Keeping the notation of the lemma we set $n(k, \ell + 1) = N(k, k, \ell)$. Given any surjective function $\rho : \Delta = [1, n(k, \ell + 1)] \to [1, k]$, we obtain a $k$-long function $f : [0, \ell]^k \to \Delta$.

Put

$$a_r = f(\underbrace{0,\ldots,0}_{r\text{ times}},\ \underbrace{\ell,\ldots,\ell}_{(k-r)\text{ times}})$$

for all $0 \le r \le k$. Due to the pigeonhole principle we find $a_r, a_s$ such that $\rho(a_r) = \rho(a_s)$. Without loss of generality let $r < s$ and define

$$h(i) = f(\underbrace{0,\ldots,0}_{r\text{ times}},\ \underbrace{i,\ldots,i}_{(s-r)\text{ times}}\ ,\underbrace{\ell,\ldots,\ell}_{(k-s)\text{ times}}).$$

We obtain for $0 \le i < \ell$ that

$$(\rho \circ h)(i) = (\rho \circ f)(\underbrace{0,\ldots,0}_{r\text{ times}},\ \underbrace{i,\ldots,i}_{(s-r)\text{ times}}\ ,\underbrace{\ell,\ldots,\ell}_{(k-s)\text{ times}}) =$$

$$(\rho \circ f)(\underbrace{0,\ldots,0}_{s\text{ times}},\ \underbrace{\ell,\ldots,\ell}_{(k-s)\text{ times}}) = (\rho \circ h)(0),$$

as $f$ is $k$-long. Moreover

$$(\rho \circ h)(\ell) = (\rho \circ f)(\underbrace{0,\ldots,0}_{r\text{ times}},\ \underbrace{\ell,\ldots,\ell}_{(k-r)\text{ times}}) = \rho(a_r) = \rho(a_s) =$$

$$(\rho \circ f)(\underbrace{0,\ldots,0}_{s\text{ times}},\ \underbrace{\ell,\ldots,\ell}_{(k-s)\text{ times}}) = (\rho \circ h)(0).$$

Whence there exists $t \in [1,k]$ such that $h([0,\ell]) \subseteq \rho^{-1}(t)$. By definition of the $k$-long map $f$ we see

$$h(i) = f(\underbrace{0,\ldots,0}_{r\text{ times}},\ \underbrace{i,\ldots,i}_{(s-r)\text{ times}}\ ,\underbrace{\ell,\ldots,\ell}_{(k-s)\text{ times}})$$

$$= a + \sum_{\nu=1}^{r} 0 \cdot d_\nu + \sum_{\nu=r+1}^{s-r} i d_\nu + \sum_{\nu=s-r+1}^{k} \ell d_\nu = \Big(a + \sum_{\nu=s-r+1}^{k} \ell d_\nu\Big) + \Big(\sum_{\nu=r+1}^{s-r} d_\nu\Big)i.$$

Thus we have shown that $h(i)$ is an arithmetic progression of length $\ell + 1$ in $\rho^{-1}(t)$. As $\rho$, being a surjective function, induces a partition $\{\rho^{-1}(1),\ldots,\rho^{-1}(k)\}$ of $\Delta$ into $k$ subsets and $\rho$ was arbitrary, the proof is complete. $\square$

### 6.3.2 Proof of the main theorem

We introduce for a ring $R$ the auxiliary notation $\Sigma_r \subseteq R$, denoting the set of all elements being representable by a sum of exactly $r$ units.

*Proof of Theorem 6.15.* Suppose we could verify the first claim, i.e. for all $r \in \mathbb{N}$ every arithmetic progression in $\Sigma_r \subseteq \mathcal{O}$ is of finite length. This immediately evidences the second claim showing that $\mathcal{O}$ cannot be $r$-good, since otherwise $\Sigma_r = \mathcal{O}$, which does contain infinite arithmetic progressions.

Using induction we start with $r = 1$. For this purpose consider $a_j = a_0 + (j-1)d \in \mathcal{O}^* = \Sigma_1$, $d \neq 0$. Due to Dirichlet's unit theorem $\mathcal{O}^*$ is a finitely generated subgroup of $\mathbb{C}^*$ as required by Theorem 5.5. The arithmetic progression may be written in the form $a_{j+1} - a_j = d$ and the left hand side interpreted as non-degenerate solutions to equation (5.4) in Theorem 5.5. Hence we obtain a bound $M_2$ restricting the progression's length.

Assume the claim holds for an $r \in \mathbb{N}$, let $M_r^*$ denote a bound for any progression's length in $\Sigma_r$. Choose $d \neq 0$ and set

$$\Omega = \{\pm u \in \mathcal{O}^* : d = u + v, \text{ with } v \in \Sigma_s \text{ and } 1 \leq s \leq 2r+1\},$$

where we require that for $u + v$ written as sum of units no subsum vanishes[5]. This set is finite[6] again due to Theorem 5.5, thus we put $\Omega =: \{x_1, \ldots, x_T\}$. Suppose now there exists an arithmetic progression $a_j = a_0 + (j-1)d$ in $\Sigma_{r+1}$ of length $W := n(T, M_r^*+1)$.

Representing the left hand side of the equation $a_{j+1} - a_j = d \neq 0$ as sum of units we find for each $0 \leq j < W$ an $\epsilon_j \in \Omega$ appearing as summand. Hence for fixed choices of $\epsilon_j$ per pair $(a_{j+1}, a_j)$ the map

$$f : \{1, \ldots, W\} \to \{1, \ldots, T\} \text{ mapping } j \mapsto t, \text{ if } \epsilon_j = x_t,$$

is well-defined. Now $\{f^{-1}(t) : 1 \leq t \leq T\}$ is a partition of $W$ into $T$ subsets. Van der Waerden's theorem yields an arithmetic progression

$$j_i = i_0 + (i-1)h \in f^{-1}(t_0)$$

for some $1 \leq t_0 \leq T$ having length $M_r^* + 1$. Now $b_i := a_{j_i} - x_{t_0}$ is an arithmetic progression[7] in $\Sigma_r$ as by construction $x_{t_0}$ is contained in the representation of $a_{j_i}$ as sum of $r+1$ units. The length of the progression $b_i$ is $M_r^* + 1$, which contradicts the induction hypothesis.                                                                                    $\square$

## 6.4 Szemerédi's theorem and density

By using Szemerédi's theorem [10] a famous, profound generalisation of van der Waerden's theorem and one of the most powerful techniques, when it comes to arithmetic progressions of rational integers, one finds a description of the density of $n$-good elements in $\mathcal{O}$.

---

[5]cf. notation of *non-degenerate* prior to Theorem 5.1

[6]$\sum_{i=2}^{2r+2} M_i$ may be used as a bound.

[7]$b_i = a_{i_0+(i-1)h} - x_{t_0} = a_0 + (i_0 + (i-1)h - 1)d - x_{t_0} = \big(d(i_0 - h - 1) + a_0 - x_{t_0}\big) + (hd)i.$

**Theorem 6.19** (Szemerédi's theorem). *For all $\epsilon > 0$ and all $\ell \in \mathbb{N}$ there exists a positive integer $N = N(\ell, \epsilon)$ such that a subset $T$ of $\mathbb{N}$ satisfying $|T| > \epsilon N$ contains an arithmetic progression of length $\ell$.*

We use Szemerédi's theorem to give another proof of van der Waerden's theorem 6.17.

*Proof of Theorem 6.17.* Let $k, \ell \in \mathbb{N}$ be given. For $\mathbb{N} \ni n \geq k$ there exists a subset $T$ in any partition of $[1, n]$ into $k$ subsets such that $|T| \geq \frac{n}{k}$. Choose some $\epsilon > 0$, and let $N$ as in Szemerédi's theorem. Now choose the number $n$ so large that $\frac{n}{k} > \epsilon N$. As this implies $|T| > \epsilon N$ Szemerédi's theorem guarantees the existence of an arithmetic progression of length $\ell$ in $T \subseteq [1, n]$. $\qquad\square$

A restatement in terms of density runs as follows:

**Corollary 6.20.** *A subset $T \subseteq \mathbb{N}$ of positive density in $\mathbb{N}$, i.e. $\lim_{n \to \infty} \frac{|T \cap [1, n]|}{n} > 0$, contains arbitrarily long arithmetic progressions.*

*Proof.* As $\lim_{n \to \infty} \frac{|T \cap [1, n]|}{n} > 0$, we find $\epsilon > 0$, such that $\frac{|T \cap [1, n]|}{n} > \epsilon$ for all $n \in \mathbb{N}$. Hence in particular $|T \cap [1, N(\ell, \epsilon)]| > \epsilon N(\ell, \epsilon)$ for all $\ell \in \mathbb{N}$, which due to Szemerédi's theorem 6.19 implies that there exist arithmetic progressions in $T \supseteq T \cap [1, N(\ell, \epsilon)]$ of arbitrary length $\ell \in \mathbb{N}$.

$\qquad\square$

Apart from Szemerédi's original paper the proofs of Furstenberg [59] and Gowers [60] are to be highlighted.

We finish the section by stating

**Proposition 6.21** ([7, Corollary 6 and Lemma 3]). *Fix a number field $K$ and some $n \in \mathbb{N}$. For all $n \in \mathbb{N}$ the set $N_N = \{x \in \mathbb{N} | x \text{ is } k\text{-good in } K, \text{ for } k \leq n\}$ has zero density in $\mathbb{N}$.*

*Proof.* Suppose $N$ has positive density, then by Szemerédi's theorem we find arbitrarily long arithmetic progressions of positive integers with each integer being a sum of *at most* $n$ units. Let such an arithmetic progression $S = \{a + md\}_{m=1}^{M}$ be partitioned into sets $\bigcup_{i=1}^{n} S_i$, such that an $S_i$ contains all $i$-good elements of $S$. Moreover set $A_i = \{m : a + md \in S_i\}$ for $i \in \{1, \ldots, n\}$. An application van der Waerden's theorem to the set $\{1, \ldots, M\} = \bigcup_{i=1}^{n} A_i$ with $M$ sufficiently large, demonstrates that we can find an arithmetic progression in one of the $A_i$'s and hence within one of the $S_i$'s. Its length increases with $M$, which in turn can be chosen arbitrarily large - a contradiction to Theorem 6.15. $\qquad\square$

# Abstract

The unit sum number $u(S)$ of a ring $S$ is defined as

$$u(S) = \begin{cases} k & S \text{ is } k\text{-good, but not } j\text{-good for all } j < k \text{ with } j, k \in \mathbb{N} \\ \omega & S \text{ is not } k\text{-good for any } k \in \mathbb{N}, \text{ but every element is a finite sum of units} \\ \infty & \text{there exists an element in } S \text{ not expressible as a finite sum of units in S} \end{cases},$$

where we say the ring $S$ is $k$-good, if every element in $S$ can be written as a sum of exactly $k$ units in S.

The thesis deals with the major results regarding specific classes of rings aiming to determine their unit sum number. It is proved that for matrix rings the unit sum number does not exceed three. The case of non-commutative, semilocal rings is completely treated by virtue of the Artin-Wedderburn structure theorem. With respect to Dedekind domains $\mathfrak{O}$, we establish a deep result by Vámos and Wiegand about an astonishing connection between the unit sum number of matrix rings over $\mathfrak{O}$ and the class number of $\mathfrak{O}$. Hereafter, an account of Jarden and Narkiewicz's recent result about algebraic rings of integers is given.

# Zusammenfassung

Die Einheitensummenzahl $u(S)$ eines Ringes $S$ ist definiert als

$$u(S) = \begin{cases} k & S \text{ ist } k\text{-gut, aber nicht } j\text{-gut für alle } j < k \text{ mit} j, k \in \mathbb{N} \\ \omega & S \text{ ist nicht } k\text{-good für irgendein } k \in \mathbb{N}, \\ & \text{aber jedes Element ist endliche Summe von Einheiten} \\ \infty & \text{es gibt ein Element aus } S, \text{ das nicht endliche Summe von Einheiten ist} \end{cases},$$

wobei der Ring $S$ $k$-gut genannt wird, falls jedes Element in $S$ als Summe von genau $k$ Einheiten in $S$ geschrieben werden kann.

Die vorliegende Arbeit beschäftigt sich mit den Hauptresultaten hinsichtlich gewisser Klassen von Ringen in Bezug auf deren Einheitensummenzahl. Es wird gezeigt, dass Matrizen eine Einheitensummenzahl $\leq 3$ aufweisen. Der Fall nicht-kommutativer, semilokaler Ringe wird unter Zuhilfenahme des Struktursatzes von Artin-Wedderburn vollständig behandelt. Bezüglich Dedekindringen $\mathfrak{O}$ werden wir ein tiefliegendes Resultat von Vámos and Wiegand über einen überraschenden Zusammenhang zwischen der Einheitensummenzahl von Matrizenringen über $\mathfrak{O}$ und der Klassenzahl von $\mathfrak{O}$ herstellen. Danach werden wir eine Abhandlung betreffs kürzlich erschienener Ergebnisse von Jarden und Narkiewicz über Ganzheitsringe bearbeiten.

# Danksagung

Wenn auch der Zeitraum zwischen Beginn der Beschäftigung mit dem vorliegenden Thema und Abgabe der Masterarbeit sich recht genau auf ein Jahr beläuft, ist mir erst in den letzten Wochen vor dem Einreichen derselbigen, dank äußerst kurzfristigen Einsatzes von Joachim Mahnkopf, eine probate Betreuung zu Teil geworden. Ihm alleine ist es geschuldet, dass die Arbeit akribisch durchbesprochen wurde - und es entbehrt wohlweislich jeglicher Übertreibung, zu bemerken, dass er der betreuerischen Aufgabe, die sich gewöhnlicherweise auf viele Monate erstreckt, zwangsweise komprimiert, aber ohne noch im geringsten an Sorgfalt einzubüßen, innerhalb kürzester Zeit verständnisvoll und engagiert nachkam. Eingedenk der glücklichen, zeitgerechten Beendigung der unter vorhergehenden, ungünstigen Umständen begonnenen Masterarbeit, bin ich also Joachim Mahnkopf zu größtem Dank verpflichtet.

Weiters gilt es dem Studienprogrammleiter Günther Hörmann meinen Dank auszusprechen, der trotz der temporalen Limitationen sehr zuvorkommend in allen administrativen Angelegenheiten agierte und auch die schlussendliche Begleitung durch Joachim Mahnkopf begünstigte.

In Dankbarkeit bin ich auch meinem guten Freund und Kommilitonen Michael Kretschy verbunden, der beim finalen Durchsehen der Arbeit mit beeindruckender Genauigkeit jeden noch so kleinen Tippfehler und jede notationelle Schwäche, sei sie von offensichtlichster oder kaum wahrnehmbarer Gestalt, eruierte und monierte, der mir, was von noch weitaus tiefgehenderer Bedeutung ist, bei allen im Laufe der Beschäftiung mit dem Thema exogen aufgetretenen Komplikationen moralische Stütze war - immer wohlwollend bemüht meine Freude an der Mathematik aufrecht zu erhalten, ja zu fördern.

Weder Bachelor- noch Masterstudium wären ohne die monetäre, niemals in Frage stehende Zuwendung durch meinen Vater möglich gewesen, ich danke ihm für seine bedingungslose Unterstützung.

# Curriculum Vitae

## Personal

| | |
|---|---|
| Name: | Thomas Blank |
| Nationality: | Austrian |

## Education

| | |
|---|---|
| 2003-2007 | Humanistisches Gymnasium Babenbergerring Wiener Neustadt |
| 2007-2011 | Bachelor Studies in Mathematics at University of Vienna |
| 2011-2014 | Master Studies in Mathematics at University of Vienna |
| | with Erasmus stay at UAB (Barcelona) |

## Employment

| | |
|---|---|
| 2013 | Tutor at the Department of Mathematics |

# Bibliography

[1] Daniel Zelinsky. Every linear transformation is a sum of nonsingular ones. *Proc. Amer. Math. Soc.*, 5:627–630, 1954.

[2] Melvin Henriksen. Two classes of rings generated by their units. *J. Algebra*, 31:182–193, 1974.

[3] Nahid Ashrafi and Peter Vámos. On the unit sum number of some rings. *Q. J. Math.*, 56(1):1–12, 2005.

[4] Lawrence S. Levy. Almost diagonal matrices over Dedekind domains. *Math. Z.*, 124:89–99, 1972.

[5] Ernst Steinitz. Rechteckige Systeme und Moduln in algebraischen Zahlkörpern. II. *Math. Ann.*, 72(3):297–345, 1912.

[6] Wolfgang Krull. Matrizen, Moduln und verallgemeinerte Abelsche Gruppen im Bereich der ganzen algebraischen Zahlen. (Beitr. z. Algebra Nr. 19.). *Sitzungsber. Heidelberger Akad. Wiss.*, 1932(2):13–38, 1932.

[7] Moshe Jarden and Władysław Narkiewicz. On sums of units. *Monatsh. Math.*, 150(4):327–332, 2007.

[8] J.-H. Evertse, H. P. Schlickewei, and W. M. Schmidt. Linear equations in variables which lie in a multiplicative group. *Ann. of Math. (2)*, 155(3):807–836, 2002.

[9] R. Swan. Van der waerdens theorem on arithmetic progressions. *Unpublished?*

[10] E. Szemerédi. On sets of integers containing no $k$ elements in arithmetic progression. *Acta Arith.*, 27:199–245, 1975. Collection of articles in memory of Juriĭ Vladimirovič Linnik.

[11] I. Martin Isaacs. *Algebra: a graduate course*, volume 100 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 2009. Reprint of

the 1994 original.

[12] Peter Vámos. 2-good rings. *Q. J. Math.*, 56(3):417–430, 2005.

[13] David Singmaster and D. M. Bloom. Rings of order four. *The American Mathematical Monthly*, 71(8):918–920, 1964.

[14] R. Raphael. Rings which are generated by their units. *J. Algebra*, 28:199–205, 1974.

[15] Peter Vámos and Sylvia Wiegand. Block diagonalization and 2-unit sums of matrices over Prüfer domains. *Trans. Amer. Math. Soc.*, 363(9):4997–5020, 2011.

[16] Irving Kaplansky. Elementary divisors and modules. *Trans. Amer. Math. Soc.*, 66:464–491, 1949.

[17] T. Y. Lam. *Serre's problem on projective modules.* Springer Monographs in Mathematics. Springer-Verlag, Berlin, 2006.

[18] Leonard Gillman and Melvin Henriksen. Rings of continuous functions in which every finitely generated ideal is principal. *Trans. Amer. Math. Soc.*, 82:366–391, 1956.

[19] Dino Lorenzini. Elementary divisor domains and Bézout domains. *J. Algebra*, 371:609–619, 2012.

[20] Morris Newman. *Integral matrices.* Academic Press, New York, 1972. Pure and Applied Mathematics, Vol. 45.

[21] Leonard Gillman and Melvin Henriksen. Some remarks about elementary divisor rings. *Trans. Amer. Math. Soc.*, 82:362–365, 1956.

[22] Moshe Roitman. The Kaplansky condition and rings of almost stable range 1. *Proc. Amer. Math. Soc.*, 141(9):3013–3018, 2013.

[23] Max D. Larsen, William J. Lewis, and Thomas S. Shores. Elementary divisor rings and finitely presented modules. *Trans. Amer. Math. Soc.*, 187:231–248, 1974.

[24] C.C. Mac Duffee. *The theory of matrices.* Chelsea publishing company, 1946.

[25] V. M. Prokip. Reduction of a set of matrices over a principal ideal domain to the Smith normal forms by means of the same one-sided transformations. In *Matrix methods: theory, algorithms and applications*, pages 166–174. World Sci. Publ., Hackensack, NJ, 2010.

[26] Olaf Helmer. The elementary divisor theorem for certain rings without chain condition. *Bull. Amer. Math. Soc.*, 49:225–236, 1943.

[27] Henri Cohen. Hermite and Smith normal form algorithms over Dedekind domains. *Math. Comp.*, 65(216):1681–1699, 1996.

[28] Robert M. Guralnick, Lawrence S. Levy, and Charles Odenthal. Elementary

divisor theorem for noncommutative PIDs. *Proc. Amer. Math. Soc.*, 103(4):1003–1011, 1988.

[29] P. B. Bhattacharya, S. K. Jain, and S. R. Nagpaul. *Basic abstract algebra.* Cambridge University Press, Cambridge, second edition, 1994.

[30] T. Y. Lam. *A first course in noncommutative rings*, volume 131 of *Graduate Texts in Mathematics.* Springer-Verlag, New York, 1991.

[31] J. Mahnkopf. Lecturenotes: Darstellungstheorie. 2011.

[32] Jens Carsten Jantzen and Joachim Schwermer. *Algebra.* Springer-Lehrbuch. Berlin: Springer. 335 p. EUR 24.95 , 2006.

[33] J. H. MacLagan Wedderburn. On Hypercomplex Numbers. *Proc. London Math. Soc.*, S2-6(1):77, 1907.

[34] Thomas W. Hungerford. *Algebra*, volume 73 of *Graduate Texts in Mathematics.* Springer-Verlag, New York, 1980. Reprint of the 1974 original.

[35] Władysław Narkiewicz. *Elementary and analytic theory of algebraic numbers.* Springer Monographs in Mathematics. Springer-Verlag, Berlin, third edition, 2004.

[36] Lawrence S. Levy and J. Chris Robson. Matrices and pairs of modules. *J. Algebra*, 29:427–454, 1974.

[37] Wolfgang M. Schmidt. Norm form equations. *Ann. of Math. (2)*, 96:526–551, 1972.

[38] Wolfgang M. Schmidt. The subspace theorem in Diophantine approximations. *Compositio Math.*, 69(2):121–173, 1989.

[39] J.-H. Evertse and R. G. Ferretti. A further improvement of the quantitative subspace theorem. *Ann. of Math. (2)*, 177(2):513–590, 2013.

[40] Gerd Faltings and Gisbert Wüstholz. Diophantine approximations on projective spaces. *Invent. Math.*, 116(1-3):109–138, 1994.

[41] H. P. Schlickewei and W. M. Schmidt. The number of solutions of polynomial-exponential equations. *Compositio Math.*, 120(2):193–225, 2000.

[42] P. Corvaja, W. M. Schmidt, and U. Zannier. The Diophantine equation $\alpha_1^{x_1} \cdots \alpha_n^{x_n} = f(x_1, \ldots, x_n)$. II. *Trans. Amer. Math. Soc.*, 362(4):2115–2123, 2010.

[43] Wolfgang M. Schmidt. The zero multiplicity of linear recurrence sequences. *Acta Math.*, 182(2):243–282, 1999.

[44] Christer Lech. A note on recurring series. *Ark. Mat.*, 2:417–421, 1953.

[45] Francesco Amoroso and Evelina Viada. On the zeros of linear recurrence sequences. *Acta Arith.*, 147(4):387–396, 2011.

[46] J.-H. Evertse and K. Győry. On the numbers of solutions of weighted unit

equations. *Compositio Math.*, 66(3):329–354, 1988.

[47] Paul Belcher. Integers expressible as sums of distinct units. *Bull. London Math. Soc.*, 6:66–68, 1974.

[48] Robert F. Tichy and Volker Ziegler. Units generating the ring of integers of complex cubic fields. *Colloq. Math.*, 109(1):71–83, 2007.

[49] Alan Filipin, Robert Tichy, and Volker Ziegler. The additive unit structure of pure quartic complex fields. *Funct. Approx. Comment. Math.*, 39(part 1):113–131, 2008.

[50] Kiran S. Kedlaya. A construction of polynomials with squarefree discriminants. *Proc. Amer. Math. Soc.*, 140(9):3025–3033, 2012.

[51] A. Y. Khinchin. *Three pearls of number theory.* Dover Publications Inc., Mineola, NY, 1998. Translated from the Russian by F. Bagemihl, H. Komm, and W. Seidel, Reprint of the 1952 translation.

[52] W. T. Gowers. A new proof of Szemerédi's theorem. *Geom. Funct. Anal.*, 11(3):465–588, 2001.

[53] E. R. Berlekamp. A construction for partitions which avoid long arithmetic progressions. *Canad. Math. Bull.*, 11:409–414, 1968.

[54] V. Chvátal. Some unknown van der Waerden numbers. In *Combinatorial Structures and their Applications (Proc. Calgary Internat. Conf., Calgary, Alta., 1969)*, pages 31–33. Gordon and Breach, New York, 1970.

[55] R. S. Stevens and R. Shantaram. Computer-generated van der Waerden partitions. *Math. Comp.*, 32(142):635–636, 1978.

[56] Michal Kouril and Jerome L. Paul. The van der Waerden number $W(2,6)$ is 1132. *Experiment. Math.*, 17(1):53–61, 2008.

[57] Michal Kouril. Computing the van der waerden number w(3, 4) = 293. *Unpublished?*, 2012.

[58] Michael D. Beeler and Patrick E. O'Neil. Some new van der Waerden numbers. *Discrete Math.*, 28(2):135–146, 1979.

[59] Harry Furstenberg. Ergodic behavior of diagonal measures and a theorem of Szemerédi on arithmetic progressions. *J. Analyse Math.*, 31:204–256, 1977.

[60] W. T. Gowers. A new proof of Szemerédi's theorem. *Geom. Funct. Anal.*, 11(3):465–588, 2001.