



SOME LATTICE POINT
PROBLEMS

by
Lewis Low M.Sc. (Sydney)

A Thesis submitted for the Degree of
Doctor of Philosophy
in The University of Adelaide,
Department of Pure Mathematics,
September, 1978.

Awarded November 1979

CONTENTS

Summary	(i)
Statement	(iv)
Acknowledgements	(v)
<u>GENERAL INTRODUCTION :</u>	1
<u>CHAPTER 1. GEOMETRICAL PRELIMINARIES</u>	
§1A : Convex sets and cones.	3
§1B : Barycentric subdivision of cones.	6
§1C : A combinatorial result.	13
<u>CHAPTER 2. AUXILIARY RESULTS ON RESIDUES</u>	15
<u>CHAPTER 3. SUBDIVISION OF n-DIMENSIONAL LATTICE CONES</u>	
§3A : Introduction to subdivision.	24
§3B : Outer index of an L-cone, and special L-cones.	31
§3C : Basic subdivisions of a complex of L-cones.	37
§3D : Minimal special subdivisions, and inner index of a complex of L-cones.	43
§3E : Results on indecomposable points.	47
<u>CHAPTER 4. SUBDIVISION OF 2-DIMENSIONAL LATTICE CONES</u>	
§4A : Preliminary.	50
§4B : The minimal basic subdivision of a 2-dimensional L-cone.	56
§4C : Basic subdivisions of a 2-dimensional L-cone.	58
§4D : Farey sequences, convergents and best approximations.	61
§4E : Note on a lemma of Schinzel.	68
Appendix to Chapter 4:	70
<u>CHAPTER 5. SUBDIVISION OF 3-DIMENSIONAL LATTICE CONES</u>	
§5A : Introduction.	73
§5B : Characterisation of 3-dimensional special L-cones.	75
§5C : Subdivision of 3-dimensional special L-cones, and complexes.	83

§5D :	Subdivision of 3-dimensional L-cones, and complexes.	90
§5E :	Applications of §5B.	91
CHAPTER 6. <u>QUADRATIC FAREY SEQUENCES</u>		
§6A :	The quadratic Farey sequence K_n .	93
§6B :	Preliminary results for K_n : elements in $(r/s, r'/s')$.	96
§6C :	General results for 2 or 3 consecutive members of K_n .	99
§6D :	Elements of K_n in $[0/1, 1/n]$.	102
§6E :	Elements of K_n in $\left[\frac{n-1}{n}, \frac{1}{1}\right]$.	103
§6F :	The smallest element of K_n in $(r/s, r'/s')$.	109
§6G :	The second smallest element of K_n in $(r/s, r'/s')$.	117
§6H :	The quadratics of form 1 at r/s .	120
CHAPTER 7. <u>AN INDEX PROBLEM OF CASSELS.</u>		
§7A :	Introduction.	124
§7B :	The set of minimal points of a 3-dimensional lattice, notation, examples.	125
§7C :	Minimal points of a 2-dimensional lattice.	131
§7D :	Distance of lattice planes from the origin.	135
§7E :	Covering domains.	139
§7F :	Bounds for the index.	144
BIBLIOGRAPHY :		149

STATEMENT

This thesis contains no material which has been accepted for the award of any other degree or diploma in any university and, to the best of my knowledge and belief, contains no material previously published or written by another person, except when due reference is made in the thesis.

ACKNOWLEDGEMENTS

I would like to thank Dr. Jane Pitman, who suggested the various topics of research, and showed a detailed interest in my work. I am indebted to her for her helpful criticisms and suggestions.

I would also like to thank A. Hesterman for information about his calculations in connection with a conjecture on quadratic Farey sequences.

Finally, I would like to thank Katherine Halsey for typing this thesis.



GENERAL INTRODUCTION.

1. The scope of the thesis.

The material in this thesis falls into three parts which are linked by their use of ideas from the geometry of numbers and their emphasis on 3-dimensional lattices.

The first part of the thesis is in the first five chapters and concerns problems on subdivision of lattice cones related to a lattice point problem of Schinzel. (A lattice cone is a cone with vertex at the origin and edges determined by lattice points.) After establishing the necessary preliminary geometrical and arithmetical results in Chapters 1 and 2, I will introduce the subdivision problems at the beginning of Chapter 3. These problems will be considered for lattice cones in n dimensions, 2 dimensions and 3 dimensions in Chapters 3, 4 and 5, respectively. Applications of the 2-dimensional results to another problem of Schinzel and of the 3-dimensional results to lattice octahedra will also be given.

The second part of this thesis concerns work arising from a problem of Mahler on the quadratic analogue of the classical Farey sequences, and my results are presented in Chapter 6. To prove one of the results there (Theorem 6C.3) another application is made of one of the main results on 3-dimensional lattices in Chapter 5.

The third part of this thesis concerns an index problem of Cassels. Our investigations leads to bounds on the index of the sublattice of a 3-dimensional lattice generated by points with minimal $|x_1| + |x_2| + |x_3|$ in each octant. The results comprise Chapter 7.

2. Preliminaries and notation. In this thesis, references are given by a number in square brackets, and refer to the bibliography at the end of the thesis. The numbering on displayed formulae, equations, and so on begins afresh with each section of a chapter. In each section of a chapter the theorems, corollaries, lemmas and propositions, and notes where the numbering includes the section number, will be numbered in one sequence. For example, in section §7D of Chapter 7, Lemma 7D.2 is followed by Theorem 7D.3, and in section §3A of Chapter 3, Note 3A.1 is followed by Lemma 3A.2. In each section of a chapter, definitions, where numbered, will be numbered in a separate sequence to that of the theorems, and so on.

Throughout this thesis, \mathbb{Z} will denote the set of all integers. Some notation of a geometrical or topological nature such as $\text{lin } A$, $\text{relint } A$ is given in §1A. The notation $\ell(a_1, \dots, a_n)$ for a lattice with basis a_1, \dots, a_n is as mentioned in §3B. Most of the definitions and notation needed for Chapters 3 to 5 appear in §3A or §3B which introduce the problem of subdivision of cones.

CHAPTER 1

GEOMETRICAL PRELIMINARIES

This chapter is on geometrical preliminaries. In §1A we discuss background material on convex sets and pointed polyhedral cones. In §1B we introduce a major tool for our later study of subdivision of lattice cones, namely, barycentric subdivision. In §1C we discuss an Euler-type combinatorial result which has applications to Chapter 5.

§1A. Convex sets and cones.

In this section we will discuss the background material on convex sets and pointed polyhedral cones which is assumed in the rest of this chapter, and in Chapters 3 to 5. We will give some lemmas for later use.

As far as possible, we will, in this thesis, use for this material, the vocabulary and definitions of Grünbaum [5], Chapter 2 and Sections 3.1, 3.2. Thus we will assume the definition and basic properties given there for the following:

Convex set, polyhedral set, polytope, face, facet, vertex, edge, extreme point, pointed cone with apex a (and hence with vertex a).

We will say that a hyperplane H in R^d cuts a set A in R^d if each of the open half-spaces determined by H contains a point of A . (Grünbaum [5], §2.2.)

We will also use the following notation of Grünbaum [5], where A is a subset of d -dimensional real space R^d :

$\text{lin } A$: the linear subspace spanned by A ,
 $\text{aff } A$: the affine subspace spanned by A ,
 $\text{dim } A$: the dimension of $\text{aff } A$, that is, one less than the maximum number of affinely independent points of A ,
 $\text{cone } A_b$: the cone with apex b generated by A ,
 $\text{conv } A$: the convex hull of A ,
 $\text{relint } A$: the relative interior of A (relative to $\text{aff } A$),
 $\text{relbd } A$: the relative boundary of A (relative to $\text{aff } A$).

We will say that two sets A, B in R^d are non-overlapping if their relative interiors are disjoint, that is if

$$\text{relint } A \cap \text{relint } B = \emptyset.$$

We will also say that a collection S of sets in R^d is non-overlapping, or that S has non-overlapping members, if each pair of distinct sets in S are non-overlapping.

Lemma 1A.1. Let A, B be convex sets in R^d such that A and B are non-overlapping, and $\text{aff } A \subseteq \text{aff } B$. Then

$$A \cap \text{relint } B = \emptyset.$$

Proof. Suppose to the contrary that $A \cap \text{relint } B \neq \emptyset$, and let $\underline{x} \in A \cap \text{relint } B$. Since $\underline{x} \in \text{relint } B$, there is an open set N of R^d such that

$$\underline{x} \in N \cap \text{aff } B \subseteq \text{relint } B.$$

Since A is convex, and $\underline{x} \in N \cap A$, there is a point \underline{y} in $N \cap \text{relint } A$ (and in fact A , $N \cap A$ and $N \cap \text{relint } A$ all have the same dimension). Hence $\underline{y} \in N \cap \text{aff } A \subseteq N \cap \text{aff } B \subseteq \text{relint } B$. Thus $\underline{y} \in \text{relint } A \cap \text{relint } B$, so A and B overlap, a contradiction.

Lemma 1A.2. Let A, B be distinct convex sets in R^d such that $A \cap B$ is a face of each of A and B . Then A, B are non-overlapping.

Proof. Suppose to the contrary that there is a point \underline{x} in $\text{relint } A \cap \text{relint } B$. Let $F = A \cap B$. Then $\underline{x} \in F \cap \text{relint } B$. Since F is a face of B , it follows that $F = B$. Similarly, $F = A$, so $A = B$, a contradiction.

The pointed cones that we will be concerned with are pointed cones with apex $\underline{0}$, and by a pointed polyhedral cone we will mean a polyhedral set which is a pointed cone with apex $\underline{0}$. Let C be a pointed polyhedral cone in R^d . Then for some $m \times d$ matrix A of rank d ,

$$C = \{\underline{x} \mid A\underline{x} \leq \underline{0}\}.$$

Also, if C has r edges and $\underline{a}_1, \dots, \underline{a}_r$ are non-zero points on distinct edges, then

$$\begin{aligned} C &= \{\lambda_1 \underline{a}_1 + \dots + \lambda_r \underline{a}_r \mid \lambda_i \geq 0 \quad \forall i\} \\ &= \text{cone}_{\underline{0}}(\underline{a}_1, \dots, \underline{a}_r), \end{aligned}$$

and $\dim C \leq r$.

Definition 1A.1. A pointed polyhedral cone C for which $\dim C$ is equal to the number of edges of C will be called a simplicial cone.

We comment that if P is an n -dimensional polytope in R^d and $\underline{b} \notin \text{aff } P$, then $C = \text{cone}_{\underline{b}}(P)$ is an $(n+1)$ -

dimensional pointed polyhedral cone, and every pointed polyhedral cone in R^d is of this form. Apart from the empty face of C , there is a one-one correspondence between the faces F of P and the faces $\text{cone}_b(F)$ of C . The cone C is a simplicial cone if and only if P is a simplex.

We complete this section by stating two lemmas which will be used in the next section. If $P \subseteq R^d$, we will write $\text{cone}_o(a, P)$ to mean $\text{cone}_o(\{a\} \cup P)$.

Lemma 1A.3. If P is an n -dimensional pointed polyhedral cone with apex b in R^d , and a is a point in R^d such that $a \notin \text{aff } P$, then $\text{cone}_b(a, P)$ is an $(n+1)$ -dimensional pointed polyhedral cone with apex b , and the faces of $\text{cone}_b(a, P)$ are the $\text{cone}_b(a, F)$, where F is a face of P , together with the faces of P .

The lemma is easily proved.

Lemma 1A.4. (Grünbaum [5], Theorem 2.6.1). If P is a polyhedral set, then a face of a face of P is again a face of P .

§1B. Barycentric subdivision of cones.

In this section we will introduce the main geometric tool to be used in Chapter 3 for subdivision of lattice cones, namely barycentric subdivision. Since it turns out that this method leads to complexes of cones, we consider subdivision of complexes of cones into complexes of subcones.

Definition. A finite collection C of polyhedral sets in R^d is a complex (in R^d) if any two members of C meet in a face of each of them.

Sometimes "complex" is defined with the additional requirement that any face of a member is again a member. (See, for example, Grünbaum [5], §3.2.) The next lemma, (which is also needed later in this section) shows that the difference is not essential.

Definition. If C is a complex, a face of C is a face of a member of C .

Lemma 1B.1. Let C be a complex in R^d . Then any two faces of C meet in a face of each of them.

Proof. Let $C_1, C_2 \in C$, and let F_1, F_2 be faces of C_1, C_2 respectively. By symmetry, it is sufficient to prove that $F_1 \cap F_2$ is a face of F_1 . By Lemma 1A.4 it is sufficient to prove that

(i) $F_1 \cap F_2$ is a face of $F_1 \cap C_2$, and

(ii) $F_1 \cap C_2$ is a face of F_1 .

To prove (i), we may suppose that $F_2 = C_2 \cap H$, where H is a hyperplane which does not cut C_2 , and hence does not cut $F_1 \cap C_2$. Thus $F_1 \cap F_2 = (F_1 \cap C_2) \cap H$ is a face of $F_1 \cap C_2$. To prove (ii), $F_1 \cap C_2 = (F_1 \cap C_1) \cap C_2 = F_1 \cap (C_1 \cap C_2) = F_1 \cap (C_1 \cap K)$, where K is R^d or is a hyperplane which does not cut C_1 and hence does not cut F_1 . Thus $F_1 \cap C_2 = F_1 \cap K$ is a face of F_1 as required.

Definition 1B.1. Let C be a complex of polyhedral sets in R^d . Then a collection \mathcal{D} of polyhedral sets is a subdivision of C if

- (a) $UC = UD$,
- (b) each member of C is a union of members of D ,
- (c) D is a complex.

The next theorem introduces barycentric subdivision, which is an important way of constructing subdivisions. We will formulate the idea for pointed polyhedral cones with apex 0 , for this is what we will need in our later investigation of lattice cones. (An analogous theory on "triangulation of polytopes" may be inferred from the correspondence between polytopes and cones that was mentioned in §1A.)

Notation. If S is a collection of sets in R^d , and $\underline{a} \in R^d$, then $\text{cone}_0(\underline{a}, S)$ denotes the set of all $\text{cone}_0(\underline{a}, A)$ with A in S .

Theorem 1B.2. (On barycentric subdivision.) Let S be a complex of pointed polyhedral cones in R^d with apex 0 , and let $\underline{a} \in \cup S$, $\underline{a} \neq 0$. Let $S(\underline{a})$ be the collection of cones obtained from S by replacing each cone C in S containing \underline{a} by all the cones $\text{cone}_0(\underline{a}, F)$, where F is a facet of C not containing \underline{a} . (Notice that if $\underline{a} \in C$, then $\underline{a} \notin F \Leftrightarrow \underline{a} \notin \text{aff } F$.) Then $S(\underline{a})$ is a subdivision of S into pointed polyhedral cones with apex 0 ; and if S consists of n -dimensional cones, so does $S(\underline{a})$.

Proof. We need only prove that $S(\underline{a})$ is a subdivision of S , as the rest follows from Lemma 1A.3. For C in S such that $\underline{a} \in C$, let us write

$$F_C = \{F \mid F \text{ is a facet of } C \text{ such that } \underline{a} \notin F\},$$

and write

$$S_1 = \bigcup_C F_C, \quad S_2 = \{C \in S \mid \underline{a} \notin C\}.$$

Then

$$S(\underline{a}) = S_2 \cup \text{cone}_0(\underline{a}, S_1).$$

We are required to prove (a), (b), (c) in Definition 1B.1 with S for C , and $S(\underline{a})$ for \mathcal{D} . It is obvious that $\cup S(\underline{a}) \subseteq S$. Hence (a) will follow from (b).

Proof of (b) : Let $C \in S$. If $\underline{a} \notin C$, then $C \in S(\underline{a})$, so suppose that $\underline{a} \in C$. It suffices to prove that

$$(1) \quad C = \cup \text{cone}_0(\underline{a}, F_C).$$

We need only prove that the left hand side is included in the right hand side, as the reverse inclusion is trivial.

Firstly, we prove that the ray $\underline{o}\underline{a}$ in C belongs to the right hand side of (1). Since $\underline{o}\underline{a}$ is contained in each cone of the union, we need only prove that $F_C \neq \emptyset$. If $C = \underline{o}\underline{a}$, then $\underline{o} \in F_C$, so $F_C \neq \emptyset$. On the other hand, if $C \neq \underline{o}\underline{a}$, then the facets of C do not all pass through $\underline{o}\underline{a}$, as \underline{o} is an extreme point of C . Hence some facet does not contain \underline{a} , and again $F_C \neq \emptyset$.

Secondly, suppose that $\underline{y} \in C$, $\underline{y} \notin \underline{o}\underline{a}$. Then $C \neq \underline{o}\underline{a}$ and, as we have just seen, the facets of C do not all pass through $\underline{o}\underline{a}$, and hence do not all contain the triangle $\underline{o}\underline{a}\underline{y}$. Rotate a ray from $\underline{o}\underline{a}$ through $\underline{o}\underline{y}$ till it meets a facet F of C not containing $\underline{o}\underline{a}\underline{y}$. Then $F \in F_C$, and $\underline{y} \in \text{cone}_0(\underline{a}, F)$. This completes the proof of (1), and of (b).

Proof of (c) : Let $D_1, D_2 \in S(\underline{a})$. We must prove that

$$(2) \quad D_1 \cap D_2 \text{ is a face of each of } D_1 \text{ and } D_2.$$

If D_1, D_2 do not contain \underline{a} , then $D_1, D_2 \in S$, so (2) follows from the fact that S is a complex. Thus we may

suppose that for some $i = 1$ or 2 , $\underline{a} \in D_i$, and

$$(3) \begin{cases} D_i = \text{cone}_O(\underline{a}, F_i) \subseteq C_i \in S, \\ F_i \in F_{C_i}, \quad \underline{a} \notin F_i, \quad \underline{a} \in C_i. \end{cases}$$

Without loss of generality suppose that (3) holds with $i = 1$. There are two cases to consider:

- (i) $\underline{a} \notin D_2$, in which case $D_2 = C_2 \in S$, or
- (ii) $\underline{a} \in D_2$, in which case (3) also holds for $i = 2$.

To prove (2) in case (i), it is sufficient, by Lemma 1A.3, to prove that

$$(4) \quad \text{cone}_O(\underline{a}, F_1) \cap C_2 = F_1 \cap C_2,$$

for by Lemma 1B.1, $F_1 \cap C_2$ is a face of each of F_1 and C_2 . To prove (2) in case (ii), it is sufficient by Lemma 1A.3, to prove that

$$(5) \quad \text{cone}_O(\underline{a}, F_1) \cap \text{cone}_O(\underline{a}, F_2) = \text{cone}_O(\underline{a}, F_1 \cap F_2),$$

for by Lemma 1B.1, $F_1 \cap F_2$ is a face of each of F_1 and F_2 .

Now to prove (4) in case (i), it is clearly sufficient to prove that if $\underline{x}_2 \in \dot{C}_2$, then

$$(6) \quad \underline{x}_2 \notin \text{relint cone}_O(\underline{a}, \dot{F}_1)$$

for any $\underline{x}_1 \in \dot{F}_1$. So suppose $\underline{x}_2 \in \dot{C}$, $\underline{x}_1 \in \dot{F}_1$. We will prove (6) by contradiction. Suppose (6) is false. Then $\underline{x}_2 \in C_1 \cap C_2$, and any hyperplane which contains \underline{x}_2 but does not cut C_1 must contain $\text{cone}_O(\underline{a}, \underline{x}_1)$ and hence \underline{a} . But $C_1 \cap C_2$ is a face of C_1 since S is a complex. Hence $C_1 \cap C_2$ contains $\text{cone}_O(\underline{a}, \underline{x}_1)$ and hence \underline{a} , contradicting the fact that $\underline{a} \notin C_2$. This completes the proof of (4).

Now to prove (5) in case (ii), we note that

$$\text{cone}_O(\underline{a}, F_i) = \bigcup_{\underline{x}_i \in \dot{F}_i} \text{cone}_O(\underline{a}, \underline{x}_i) \quad (i = 1, 2).$$

Hence the left hand side of (5) is equal to the union of intersections

$$\text{cone}_0(\underline{a}, \underline{x}_1) \cap \text{cone}_0(\underline{a}, \underline{x}_2),$$

where $\underline{x}_1 \in \dot{F}_1$, $\underline{x}_2 \in \dot{F}_2$. Thus it is sufficient to prove that (6) holds for $\underline{x}_1 \in \dot{F}_1$, $\underline{x}_2 \in \dot{F}_2$, and we do this by contradiction. Suppose (6) is false. Then, as in case (i), $C_1 \cap C_2$ contains $\text{cone}_0(\underline{a}, \underline{x}_1)$. Hence $\text{cone}_0(\underline{a}, \underline{x}_1) \subseteq C_2$, contradicting the assumption that the point \underline{x}_2 in $\text{relint } \text{cone}_0(\underline{a}, \underline{x}_1)$ lies on a facet F_2 of C_2 such that $\underline{a} \notin \text{aff } F_2$. This completes the proof of (5), and the proof of the theorem.

Definition 1B.2. The complex $S(\underline{a})$ in Theorem 1B.2 is said to be obtained from S by barycentric subdivision with division point \underline{a} .

Barycentric subdivision of a simplicial cone. We will describe this case in terms of generators. Let

$$C = \text{cone}_0(\underline{a}_1, \dots, \underline{a}_n)$$

be an n -dimensional simplicial cone (Definition 1A.1) with non-zero points $\underline{a}_1, \dots, \underline{a}_n$ on distinct edges, and let

$$\underline{a} = \lambda_1 \underline{a}_1 + \dots + \lambda_n \underline{a}_n \in C, \quad \underline{a} \neq \underline{o}.$$

Let

$$I = \{i \mid \lambda_i > 0\},$$

and define

$$C_i = \text{cone}_0(\underline{a}_1, \dots, \underline{a}_{i-1}, \underline{a}, \underline{a}_{i+1}, \dots, \underline{a}_n)$$

for $i \in I$. It is easily seen that barycentric subdivision of C with division point \underline{a} gives

$$C(\underline{a}) = \{C_i \mid i \in I\},$$

and that if

$$b = \mu_1 a_1 + \dots + \mu_n a_n \in C, \quad b \neq o,$$

then

$$b \in C_i \Leftrightarrow \frac{\mu_i}{\lambda_i} = \min_{j \in I} \frac{\mu_j}{\lambda_j}.$$

Definition. A complex of pointed polyhedral cones is simplicial if all its members are simplicial.

Theorem 1B.3. (Simplicial subdivision.) Let S be a complex of n -dimensional pointed polyhedral cones in R^d with apex o . Then

(i) S is simplicial if and only if $S(\underline{a}) = S$ whenever \underline{a} is a division point on an edge of S .

(ii) Some finite sequence of successive barycentric subdivisions yields a subdivision T of S into n -dimensional simplicial cones such that the edges of T are just the edges of S .

Proof. (i) It is clearly sufficient to prove the result in the case when S consists of a single cone, say $S = \{C\}$. The necessity of the condition is obvious. We will prove sufficiency by mathematical induction. The result is trivially true when $n=1$ or 2 . Suppose that the result is true when $n=k \geq 2$, and let $n=k+1$. We assume that $C(\underline{a}) = C$ whenever \underline{a} is a non-zero point on an edge of C , and we are required to prove that C is simplicial. Suppose not. Then C has more than $k+1$ edges. Choosing any \underline{a} as above, we deduce from $C(\underline{a}) = C$ that there is a facet D of C such that $\underline{a} \notin \text{aff } D$ and $C = \text{cone}_o(\underline{a}, D)$. Then D is k dimensional with more than k edges, and so is not simplicial. By induction hypothesis, there is a non-zero point \underline{b} on an edge of D

such that $D(\underline{b})$ has at least two members D_1, D_2 . Thus $C(\underline{b})$ has at least two members $\text{cone}_{\underline{a}}(\underline{a}, D_1)$, $\text{cone}_{\underline{a}}(\underline{a}, D_2)$, contradicting the fact that $C(\underline{b}) = C$.

(ii) Beginning with the complex S , make repeated use of barycentric subdivision, where the division point is on an edge of the complex. The set of edges remains unchanged. By (i), the number of members of the complex increases at each step unless a simplicial complex is reached. This eventually happens because the number of simplexes whose edges are edges of S is finite.

The theory of subdivision cannot be simplified by choosing only division points in the relative interiors of cones. For example, if P is a 3-dimensional square pyramid in R^4 such that $\underline{a} \notin \text{aff } P$, then $\text{cone}_{\underline{a}}(P)$ cannot be subdivided into simplicial cones by repeated barycentric subdivision using only division points in the relative interiors of cones.

§1C. A combinatorial result.

We present here a combinatorial result which essentially belongs to topology. We will need it in Chapter 5.

Lemma 1C.1. Let S be a complex of 3-dimensional simplicial cones in R^3 such that US is simply connected. Then

$$2E_i + E_b = N + 2,$$

where E_i is the number of edges of S in the relative interior of US , E_b is the number of edges on the relative boundary of US , and N is the number of members of S .

Proof. The result is easily proved by induction on

N. It is of course equivalent to the well-known result that the number of primitive triangles in a primitive triangulation of a simple polygon is

$$2V_i + V_b - 2,$$

where V_i, V_b are the number of vertices in the interior and boundary respectively of the polygon. (See, for example, Gaskell, Klamkin and Watson [4].)

CHAPTER 2

AUXILIARY RESULTS ON RESIDUES

In this chapter we will obtain results on residue sets which are needed later (for the proofs of Proposition 3E.1 and Lemma 5B.1). The main result is Proposition 2.6 below.

Let m be a fixed positive integer greater than 1, and for any integer x let \bar{x} denote the least non-negative residue of x modulo m , so that

$$0 \leq \bar{x} < m, \quad x \equiv \bar{x} \pmod{m}.$$

Definition 2.1. For x in \mathbb{Z} such that $1 \leq x \leq m-1$, let ρ_x be the function with domain \mathbb{Z} such that for r in \mathbb{Z} ,

$$\rho_x(r) = \begin{cases} 1 & \text{if } \overline{rx} = \overline{(r-1)x} + x - m \\ 0 & \text{if } \overline{rx} = \overline{(r-1)x} + x. \end{cases}$$

We note that for r in \mathbb{Z}

$$(1) \quad \overline{rx} - \overline{(r-1)x} = x - m\rho_x(r),$$

and the following properties of ρ_x are thus immediate.

Lemma 2.1. For x and r in \mathbb{Z} and $1 \leq x \leq m-1$ we have:

- (i) $\rho_x(r) = \rho_x(\bar{r})$;
- (ii) $\rho_x(r) = 1 \Leftrightarrow \overline{rx} < x$, and
 $\quad = 0 \Leftrightarrow \overline{rx} \geq x$;
- (iii) $\sum_{r=1}^m \rho_x(r) = x$;
- (iv) $\rho_1(r) = 0$ if $\bar{r} \neq 0$;
- (v) $\rho_x(0) = 1$;
- (vi) $\rho_x(1) = 0$;
- (vii) $\rho_x(2) = 1 \Leftrightarrow x \geq \frac{m}{2}$.

In view of (i), (v), (vi), we now consider the restriction $\rho_x|S$ of ρ_x to the subset S of Z given by

$$(2) \quad S = \{2, 3, \dots, m-1\}.$$

Lemma 2.2. Let S be as in (2), and let ρ_x be as in Definition 2.1.

- (i) $\rho_1|S = 0$ (the zero function on S).
- (ii) $\rho_{m-1}|S = 1$ (the constant function 1 on S).
- (iii) If $(a, m) = 1$, $1 \leq a \leq m-1$, then

$$(\rho_a|S) + (\rho_{m-a}|S) = 1 (= \rho_{m-1}|S).$$

- (iv) For $1 \leq x \leq m-1$,

$$\sum_{r \in S} \rho_x(r) = x-1.$$

- (v) Suppose a_1, a_2, \dots, a_k, b are integers such that $1 \leq a_i \leq m-1$ ($i = 1, \dots, k$), $1 \leq b \leq m-1$, and

$$\sum_{i=1}^k (\rho_{a_i}|S) = \rho_b|S.$$

Then

$$\sum_{i=1}^k a_i = b + k - 1.$$

- (vi) Suppose a_1, a_2, \dots, a_k are integers such that $1 \leq a_i \leq m-1$ ($i = 1, \dots, k$), and

$$\sum_{i=1}^k a_i > (k-1)m - (k-2).$$

Then for some r in S ,

$$\rho_{a_1}(r) = \dots = \rho_{a_k}(r) = 1,$$

that is,

$$\overline{ra_1} < a_1, \dots, \overline{ra_k} < a_k.$$

(Part (vi) is needed later in §3E.)

Proof. Parts (i), (ii) are immediate from the definitions, (iv) is immediate from Lemma 2.1(iii), (v), (vi). Part (v) is immediate from (iv).

(iii) Since $\overline{ra} + \overline{r(m-a)} \equiv 0 \pmod{m}$, $\overline{ra} + \overline{r(m-a)}$ equals 0 if both summands are 0, and equals m otherwise. If $(a, m) = 1$, the latter is the case for all r in S , and the result follows from Lemma 2.1(ii).

$$\begin{aligned} \text{(vi)} \quad \sum_{r \in S} (\rho_{a_1}(r) + \dots + \rho_{a_k}(r)) &= \sum_{i=1}^k a_i - k \quad \text{by (iv)} \\ &> (k-1)(m-2) \quad \text{by hypothesis,} \\ &= (k-1) \# S, \end{aligned}$$

so for some r in S ,

$$\rho_{a_1}(r) + \dots + \rho_{a_k}(r) = k.$$

Thus $\rho_{a_1}(r) = \dots = \rho_{a_k}(r) = 1$, and by Lemma 2.1(ii), $\overline{ra_1} < a_1, \dots, \overline{ra_k} < a_k$. This completes the proof of Lemma 2.2.

Let $[y]^+$ denote the least integer greater than or equal to y , that is,

$$[y]^+ - 1 < y \leq [y]^+.$$

Lemma 2.3 (i) For all x (such that $1 \leq x \leq m-1$), the t^{th} smallest positive value of r for which $\rho_x(r) = 1$ is

$$\left[\frac{tm}{x} \right]^+.$$

$$\text{(ii)} \quad \left[\frac{(t+1)m}{x} \right]^+ = \left[\frac{tm}{x} \right]^+ + \left[\frac{m}{x} \right]^+ - \varepsilon$$

where $\varepsilon = 0$ or 1 .

Proof. (i) This is immediate from the definition of ρ_x .

(ii) This is just a special case of the general result that for any real numbers y, z ,

$$[y+z]^+ = [y]^+ + [z]^+ - \varepsilon,$$

where $\varepsilon = 0$ or 1 .

Clearly the behaviour of $\rho_x|_S$, where S is as in (2), can be formulated in terms of the set

$$\begin{aligned} S_x &= \{r \mid r \in S, \rho_x(r) = 1\} \\ &= \{r \mid r \in S, \overline{rx} < x\} \quad (1 \leq x \leq m-1), \end{aligned}$$

of which $\rho_x|_S$ is the indicator. The following lemma is most easily formulated in this way, and is the essence of the proof of Proposition 2.6(ii).

Lemma 2.4. Let b, c be integers such that $1 \leq b \leq m-1$, $1 \leq c \leq m-1$. Suppose that $S_b \subseteq S_c$, $S_b \neq S_c$, and $\left[\frac{m}{b}\right]^+ = \left[\frac{m}{c}\right]^+$. Then $\left[\frac{m}{c}\right]^+ = 2$ (that is, $\rho_c(2) = 1$, or $c \geq \frac{m}{2}$).

Proof. The condition $S_b \subseteq S_c$ means that $\rho_b(r) \leq \rho_c(r)$ for all positive r . Since $S_b \neq S_c$, there is some $r = r_0$ in S such that $\rho_b(r_0) = 0$, $\rho_c(r_0) = 1$, and $\rho_b(r) = \rho_c(r)$ for r in S with $r < r_0$. If r_0 is the (t_0+1) -th smallest integer in S for which $\rho_c(r_0) = 1$, then by Lemma 2.3(i),

$$\left[\frac{tm}{b}\right]^+ = \left[\frac{tm}{c}\right]^+ \quad \text{for } t = 1, \dots, t_0,$$

and

$$\left[\frac{(t_0+1)m}{b}\right]^+ = \left[\frac{(t_0+1)m}{c}\right]^+ + \eta = r_0 + \eta = r_1 \quad \text{say} \\ (\leq m),$$

where $\eta > 0$, and by assumption, $t_0 \geq 1$. Applying Lemma 2.3(ii) with $t = t_0$, and $x = b, c$, we get

$$\begin{aligned} \left[\frac{(t_0+1)m}{b}\right]^+ &= \left[\frac{t_0 m}{b}\right]^+ + \left[\frac{m}{b}\right]^+ - \varepsilon_1, \\ \left[\frac{(t_0+1)m}{c}\right]^+ &= \left[\frac{t_0 m}{c}\right]^+ + \left[\frac{m}{c}\right]^+ - \varepsilon_2, \end{aligned}$$

where $\varepsilon_1, \varepsilon_2 = 0$ or 1 . Hence $\eta = \varepsilon_2 - \varepsilon_1 = 1$ (and $\varepsilon_2 = 1, \varepsilon_1 = 0$) and $r_1 = r_0 + 1$. Since $\rho_b(r_1) = 1$ (by Lemma 2.3(i)), $\rho_c(r_1) = 1$, and so $\left[\frac{(t_0+2)m}{c}\right]^+ = r_1 = r_0 + 1$.

But by Lemma 2.3(ii) with $t = t_0 + 1$, $x = c$,

$$r_1 = r_0 + \left\lceil \frac{m}{c} \right\rceil^+ - \varepsilon_3,$$

where $\varepsilon_3 = 0$ or 1 . Hence $\left\lceil \frac{m}{c} \right\rceil^+ = 1 + \varepsilon_3$, which must be equal to 2 since $m > c$.

We shall need one further lemma for our main proposition.

Lemma 2.5. Let S be as in (2), and let a_1, \dots, a_k be integers such that $1 \leq a_i \leq m-1$ for $i = 1, \dots, k$.

Then the statement

$$(3) \quad \sum_{i=1}^k (\rho_{a_i} | S) = 1 \quad (= \rho_{m-1} | S)$$

$$(\text{that is, } \sum_{i=1}^k \rho_{a_i}(r) = 1 \quad (= \rho_{m-1}(r)) \text{ for } r = 2, \dots, m-1)$$

is equivalent to

$$(4) \quad \sum_{i=1}^k \overline{ra_i} = m + r(k-2) \quad \text{for } r = 1, 2, \dots, m-1.$$

Proof. From equation (1) (following Definition 2.1) we obtain the identity

$$(5) \quad \sum_{i=1}^k \overline{ra_i} - \sum_{i=1}^k \overline{(r-1)a_i} = \sum_{i=1}^k a_i - m \sum_{i=1}^k \rho_{a_i}(r).$$

Suppose that (3) holds. Then Lemma 2.2(v) gives

$$(6) \quad \sum_{i=1}^k a_i = m + k - 2,$$

which is (4) with $r = 1$. Condition (4) follows for all r by induction using (5).

Conversely, suppose that (4) holds. Then (3) follows from (5) on applying (4) with r replaced by r , $r-1$ and 1 ($2 \leq r \leq m-1$).

The following is the main result of this chapter.

Proposition 2.6. Let S be as in (2), and let a_1, \dots, a_k be integers such that

$$1 \leq a_1 \leq \dots \leq a_k \leq m-1,$$

and

$$(4) \quad \sum_{i=1}^k \overline{ra_i} = m + r(k-2) \quad \text{for } r = 1, 2, \dots, m-1.$$

(i) Then

$$(a_i, m) = 1 \quad \text{for } i = 1, \dots, k.$$

(ii) Also,

$$1 = a_1 = \dots = a_{k-2} \leq a_{k-1} \leq a_k,$$

$$a_{k-1} + a_k = m,$$

and $a_{k-1} = a_k$ only if $m = 2$.

Proof. (i) Suppose that for $i = k$,

$$(7) \quad (a_k, m) = \delta > 1.$$

Let $\delta' = m/\delta$. For integral x , let $\langle x \rangle$ be given by

$$\langle x \rangle \equiv x \pmod{\delta}, \quad 0 \leq \langle x \rangle \leq \delta - 1,$$

while, as before, \bar{x} is given by

$$\bar{x} \equiv x \pmod{m}, \quad 0 \leq \bar{x} \leq m-1.$$

Let s be an integer such that $1 \leq s \leq \delta - 1$. From (7),

$$(8) \quad \overline{s\delta'a_k} = 0.$$

Also, as is easily checked,

$$(9) \quad \langle sx \rangle = \overline{s\delta'x/\delta'}.$$

From (4) with $r = s\delta'$, we obtain on using (8), dividing by δ' , and using (9), that

$$(10) \quad \sum_{i=1}^{k-1} \langle sa_i \rangle = \delta + s(k-2) \quad \text{for } 1 \leq s \leq \delta - 1.$$

Let us see if this condition is of the same form as (4).

We may of course replace each a_i in (10) by $\langle a_i \rangle$, and

$$0 \leq \langle a_i \rangle \leq \delta - 1 \quad \text{for } i = 1, \dots, k-1.$$

Also, $\langle a_k \rangle = 0$ by (7). Now if $1 \leq i \leq k$, and $\langle a_i \rangle = 0$, then $\delta | a_i, m | \delta' a_i$, so $\rho_{a_i}(\delta') = 1$. Thus, if there were more than one value of i with $1 \leq i \leq k$ for which $\langle a_i \rangle = 0$, then

$$\sum_{i=1}^k \rho_{a_i}(\delta') > 1.$$

By Lemma 2.5 (the (4) \Rightarrow (3) part), this contradicts (4).

Thus

$$1 \leq \langle a_i \rangle \leq \delta - 1 \quad \text{for } i = 1, \dots, k-1.$$

Thus, (10) is a similar condition to (4), but with δ in place of m , $k-1$ in place of k , and $k-2$ instead of the corresponding $k-3$ on the right hand side. By the corresponding argument to that of the proof in Lemma 2.5 that (4) \Rightarrow (3) \Rightarrow (6), we obtain

$$\langle a_1 \rangle + \dots + \langle a_{k-1} \rangle = \delta + (k-3),$$

which now contradicts (10) when $s=1$. This contradiction shows that (i) is true for $i=k$.

(ii) By Lemma 2.5, we have

$$(3) \quad \sum_{i=1}^k (\rho_{a_i} | S) = 1 = \rho_{m-1} | S.$$

If $a_{k-1} = 1$, then for $i \leq k-1$, $a_i = 1$, so $\rho_{a_i} | S = 0$, so $\rho_{a_k} | S = \rho_{m-1} | S$, and $a_k = m-1$ (by Lemma 2.1(iii), for example), and the result follows. Now suppose that $a_{k-1} > 1$. Then $\rho_{a_{k-1}} | S \neq 0$, and $\rho_{a_k} | S \neq \rho_{m-1} | S$. From (3) and the assumption that $a_1 \leq \dots \leq a_k$, we see that

$$(11) \quad \left[\frac{m}{a_k} \right]^+ = \left[\frac{m}{m-1} \right]^+ = 2, \quad a_k \geq \frac{m}{2}.$$

Since $(a_k, m) = 1$ by (i) for $i=k$, Lemma 2.2(iii) gives

$$\rho_{a_k} | S + \rho_{m-a_k} | S = 1.$$

Hence (3) gives

$$\sum_{i=1}^{k-1} (\rho_{a_i}|S) = \rho_{m-a_k}|S.$$

Now suppose, if possible, that a_1, \dots, a_{k-2} are not all equal to 1. Then $\rho_{a_{k-1}}|S \neq \rho_{m-a_k}|S$, and from Lemma 2.4 with $b = a_{k-1}$ and $c = m - a_k$, it follows that $a_{k-1}, m - a_k \geq m/2$, so $a_k \leq m/2$. From (11), $a_k = m/2$. As $(a_k, m) = 1$, it follows that $m = 2$. This is a contradiction, as (ii) (and (i)) are trivially true for $m = 2$. Thus a_1, \dots, a_{k-2} are all equal to 1, so (3) becomes

$$\rho_{a_{k-1}}|S + \rho_{a_k}|S = 1.$$

Hence, by Lemma 2.2(v),

$$a_{k-1} + a_k = m$$

and since $(a_k, m) = 1$, $a_{k-1} \neq a_k$ unless $m = 2$. This completes the proof of (ii).

Finally, (i) for $1 \leq i \leq k-1$ follows immediately from (i) for $i = k$ and (ii).

Note 1. In Chapter 5, we will be using Proposition

2.6 in the case when $k = 3$. If we were to prove the result in this case only, the only part of the proof which simplifies is the proof of (i): If $k = 3$, and $(a_3, m) = \delta > 1$, then $\delta' = \frac{m}{\delta}$,

$$\overline{\delta' a_1} + \overline{\delta' a_2} = m + \delta' > m$$

$$\overline{(m - \delta') a_1} + \overline{(m - \delta') a_2} = m + (m - \delta') > m,$$

contradicting the fact that the sum of the left hand sides is less than or equal to $2m$.

Note 2. Proposition 2.6 says that the only decompositions of $\rho_{m-1}|S$ into the sum of two or more $\rho_x|S$ are of the form

$$\rho_{m-1}|S = \rho_a|S + \rho_{m-a}|S + 0 + 0 + \dots$$

where $(a, m) = 1$. Since condition (3) is equivalent to

$$\sum_{i=1}^{k-1} \rho_{a_i} |S = \rho_{m-a_k} |S$$

provided $(a_k, m) = 1$, it follows that for a_k such that $(a_k, m) = 1$ and $a_k \neq m-1$, $\rho_{a_k} |S$ has only the trivial decompositions

$$\rho_{a_k} |S = \rho_{a_k} |S + 0 + 0 + \dots$$

It is an obvious question whether there is a non-trivial decomposition

$$(12) \quad \rho_{a_1} |S + \dots + \rho_{a_{k-1}} |S = \rho_{a_k} |S$$

when $(a_k, m) = \delta > 1$. For this it is clearly necessary that at most one of the (a_i, m) , $i \leq k-1$, exceeds 1. In fact we can prove by methods similar to those used in the proof of Proposition 2.6(i), that for any such decomposition (12), $a_{i_0} = h_{i_0} \delta$ for some $i_0 \neq k$, $a_i = h_i \delta + 1$ for $i \neq i_0, k$, and $a_k = (h_1 + \dots + h_{k-1}) \delta$.

CHAPTER 3

SUBDIVISION OF n -DIMENSIONAL LATTICE CONES

In §3A I will introduce the subdivision problem for lattice cones, and outline the material covered in the rest of Chapter 3, and in Chapters 4, and 5. I will begin with some necessary definitions, show how a main result of Chapter 1 may be applied to lattice cones, and then discuss a problem of Schinzel, which can be solved by subdivision of lattice cones.

§3A Introduction to subdivision.§3A.1. Preliminary definitions and results.

We will assume the geometrical preliminaries of Chapter 1 concerning complexes of polyhedral sets in R^d , particularly the results concerning subdivision of complexes of pointed polyhedral cones. We will also use basic results concerning indices of sublattices of point lattices in R^d , which will be stated in §3B.

An s -dimensional (point) lattice in R^d is defined to be the set of all linear combinations with integer coefficients, of s linearly independent points in R^d . A point a is said to be a primitive point of the lattice L if $ra \in L$ if and only if $r \in \mathbb{Z}$.

Let L be a fixed lattice in R^d .

Definition. An (n -dimensional) L-cone in R^d is an (n -dimensional) pointed polyhedral cone C (with apex 0) in R^d , such that each edge of C contains a non-zero point (and hence a unique primitive point) of L .

The simplest type of L -cone is now defined, and then a special case.

Definition. A simplicial L-cone is a simplicial cone which is also an L-cone.

Definition. A basic L-cone is an n -dimensional simplicial L-cone C in R^d such that the set of primitive points on the n edges of C are a basis for the n -dimensional lattice $L' = L \cap \text{lin } C$.

Thus, C is a basic L-cone if and only if C is a basic L' -cone. Also, the faces of a simplicial (or basic) L-cone are clearly simplicial (or basic) L-cones.

We notice that if we perform barycentric subdivision on an L-cone, using points of L as division points, the resulting subcones will be L-cones.

Note 3A.1. Henceforth we will consider only cones which are L-cones. We will be considering complexes of n -dimensional L-cones, and the only subdivisions of these that we will consider are subdivisions into n -dimensional L-cones.

From Theorem 1B.3 on subdivision of a complex of cones, we immediately have the following result.

Lemma 3A.2. Let S be a complex of n -dimensional L-cones in R^d . Then by repeated barycentric subdivision, we may obtain a subdivision S_1 into n -dimensional simplicial L-cones, such that the edges of S_1 are just the edges of S .

The significance of an n -dimensional simplicial L-cone C being basic, is explained in the following lemma, which is an immediate consequence of the definitions.

Lemma 3A.3. If C is basic, then every point of $L \cap C$ is expressible as a linear combination of the n

primitive points on the edges of C , with non-negative integral coefficients. (The converse is also true, as will be seen from Proposition 3B.2(ii) below.)

§3A.2. The subdivision problem, and Schinzel's problem.

Our investigation of subdivision of n -dimensional L -cones in \mathbb{R}^d will solve the following problems, and provide various bounds associated with them.

The Subdivision Problem: To find a basic subdivision of a given n -dimensional L -cone, that is, a subdivision into n -dimensional basic L -cones. (See Note 3A.1. above.)

Schinzel's Problem: Given an n -dimensional L -cone C in \mathbb{R}^d , to find a finite subset B of $L \cap C$ such that every point \tilde{x} in $L \cap C$ is expressible as a non-negative integral linear combination of an n -element subset $B_{\tilde{x}}$ of B .

According to W. Schmidt [12], A. Schinzel considered the special case of the latter problem, where $d = n$, L is an n -dimensional sublattice of \mathbb{Z}^n , and C is the non-negative orthant in \mathbb{R}^n , that is, the simplicial cone generated by the primitive points of L lying on the positive coordinate axes of \mathbb{R}^n . Schinzel conjectured that finite B existed in this case, and Schmidt [12] proved this conjecture by a (non-constructive) compactness argument using induction on n .

It is clear from Lemma 3A.3 that a solution of the subdivision problem will yield a solution of Schinzel's problem, since we may take B to be the set of primitive lattice points on the edges of the subdivision and $B_{\tilde{x}}$ to be the lattice basis determined by a basic cone of the

subdivision containing \underline{x} . Moreover, Lemma 3A.2 reduces the subdivision problem (and hence Schinzel's problem) for general L-cones to that for simplicial L-cones. I obtained a simple constructive solution of Schinzel's problem by this method, and a preliminary account of this work is given in Low [9].

Remark 1. If in Schinzel's problem, $B_{\underline{x}}$ is an $(n+1)$ -element subset of B , rather than an n -element subset, the problem becomes trivial. For in the case of simplicial C , to which the problem is reduced by Lemma 3A.2, if $\underline{a}_1, \dots, \underline{a}_n$ is the set of primitive lattice points on the edges of C , and $\underline{x}_1, \dots, \underline{x}_m$ are the coset representatives of the sublattice of $L \cap \text{lin } C$ generated by $\underline{a}_1, \dots, \underline{a}_n$ chosen in the form

$$\alpha_1 \underline{a}_1 + \dots + \alpha_n \underline{a}_n$$

where $0 \leq \alpha_i < 1$ for $i = 1, \dots, n$, then any \underline{x} in C is expressible in the form

$$\underline{x}_i + \lambda_1 \underline{a}_1 + \dots + \lambda_n \underline{a}_n,$$

where the λ_j are non-negative integers.

Remark 2. If ϕ is a non-singular linear transformation on R^d , then the subdivision problem and Schinzel's problem for the L-cone C are clearly equivalent to the corresponding problems for the $\phi(L)$ -cone $\phi(C)$. When C is a simplicial L-cone, this idea may be applied to put C in some sort of standard form, for example, that in which the primitive points of L on the edges of C are the natural unit vectors $\underline{e}_1, \dots, \underline{e}_n$. Although such a standard form would be convenient in giving examples of L-cones, it does not help us in solving the subdivision problem, since subcones of a cone in standard form would not themselves be in standard form.

§3A.3. Various definitions.

We shall now introduce some concepts needed to continue our discussion of Schinzel's problem and the subdivision problem.

Let C and D be finite collections of n -dimensional L -cones such that $UC = UD$. We call D a decomposition of C if each cone in C is a union of members of D . We call D a dissection of C if it is a decomposition of C into non-overlapping L -cones. (The term non-overlapping was defined in §1A.) A subdivision of a complex C (see Definition 1B.1) is thus a decomposition of C which is also a complex. By Lemma 1A.2, any subdivision of C is also a dissection of C .

Now in order to solve Schinzel's problem it is not necessary to restrict consideration to subdivisions into basic cones. Decompositions into basic cones would suffice, with the sets B, B_x being derived in the same way as before. In fact, Schmidt used decompositions, and in Low [9], for the sake of simplicity, I stated and proved the results in terms of dissections (although the word dissection was not used there), even though the method used, namely, barycentric subdivision, led naturally to subdivisions. Dissections are easier to handle than subdivisions because a dissection of C is just the union of dissections of the members of C , and so to obtain dissections, it suffices to consider barycentric subdivisions of L -cones, rather than of complexes of L -cones, even if some division point used is not in the relative interior of the relevant cone.

Definition. Let C be an L -cone and let $\dot{L} = L \sim \{0\}$. A point \underline{x} in $\dot{L} \cap C$ is decomposable in $\dot{L} \cap C$ (or simply in C) if there exist $\underline{x}_1, \underline{x}_2$ in $\dot{L} \cap C$ such that $\underline{x} = \underline{x}_1 + \underline{x}_2$. Otherwise, \underline{x} is indecomposable in $\dot{L} \cap C$ (or in C).

The notation $I(C)$ will be used for the set of indecomposable points of C .

Definition. For $\underline{a} = (a_1, a_2, \dots, a_d)$ in R^d the height $H(\underline{a})$ of \underline{a} is defined to be

$$H(\underline{a}) = \max_i |a_i|.$$

§3A.4. Plan of the rest of Chapter 3 and of Chapters 4 and 5.

In the remainder of Chapter 3 I will consider n -dimensional L -cones for general n . In §3B I will give the necessary preliminary results on the (outer) index of an L -cone, and on special L -cones. In §3C I will solve the subdivision problem and obtain crude bounds which in the case of Schinzel's problem are bounds for

- (i) the number of points in B
- (ii) the number of distinct $B_{\underline{x}}$
- (iii) the height of points in B .

In §3D we will study a refinement of the method of §3C, which would appear to lead to more economical bounds (and certainly does in the case of 2- and 3-dimensional L -cones, as we shall see later). In §3E we collect together results on indecomposable points.

Let S be a decomposition of an n -dimensional L -cone C into basic L -cones, and let $B(S)$ be the set of primitive points on the edges (of the members) of S . Then it is obvious that $B(S) \supseteq I(C)$. In Chapters 4 and 5 I will discuss 2-, and 3-dimensional L -cones, and will find in

each case a subdivision S of C into basic L-cones for which $B(S) = I(C)$. These results will yield sharp bounds for the numbers in (i), (ii) and (iii) for Schinzel's problem (in fact the best possible bounds for the case of subdivisions). I will also discuss there the equivalence of the concepts of decomposition, dissection and subdivision under certain conditions.

Note (added August, 1979).

Schinzel's problem is also considered in the (unpublished) thesis
R.A. Lee, A quantitative solution of a problem of Schinzel, Ph.D. thesis, University of Colorado (1970) (available through University Microfilms, London), which concentrates mainly on a different aspect of the problem.

§3B. Outer index of an L-cone, and special L-cones.

In this section, we discuss concepts needed for tackling the subdivision problem. Again, L will be a fixed lattice in R^d of dimension less than or equal to d . We will review the basic results on the index of an n -dimensional sublattice of an n -dimensional lattice. We will define the base points of an L-cone, the parallelotope of a simplicial L-cone, and the index of simplicial and the (outer) index of general L-cones, and give some results in terms of these ideas. The section ends with a discussion of special cones, which are relevant to §3D. All these concepts are needed for our investigation of subdivisions of L-cones in Chapters 3 to 5.

Notation for lattices. If $\underline{b}_1, \dots, \underline{b}_n$ are linearly independent points in R^d , then the n -dimensional lattice $N = \{\lambda_1 \underline{b}_1 + \dots + \lambda_n \underline{b}_n \mid \lambda_i \in \mathbb{Z} \ \forall i\}$ with basis $\underline{b}_1, \dots, \underline{b}_n$ is denoted by $N = \ell(\underline{b}_1, \dots, \underline{b}_n)$.

Index of sublattice.

Definition 3B.1. If $M = \ell(\underline{a}_1, \dots, \underline{a}_n)$ is a sublattice of $N = \ell(\underline{b}_1, \dots, \underline{b}_n)$ in R^d , then the index $(N:M)$ is the subgroup index, that is, the number of cosets of M in N .

Proposition 3B.1. Let M, N be as in Definition 3B.1, and let the integral matrix $\Lambda = (\lambda_{ij})$ be given by

$$\underline{a}_j = \sum_{i=1}^n \lambda_{ij} \underline{b}_i, \quad j = 1, \dots, n.$$

Then

$$(N:M) = |\det \Lambda|,$$

and in the particular case when $d = n$,

$$(N:M) = \frac{|\det(\underline{a}_1, \dots, \underline{a}_n)|}{|\det(\underline{b}_1, \dots, \underline{b}_n)|}$$

The result is easily seen from Cassels [3] I.2.2, even though Cassels considers only those lattices in R^d which are d -dimensional (that is, $d = n$), and in fact takes the last equation as a definition, and proves the coset result as a lemma.

In order to state the lemma connecting the index of a sublattice, and linear combinations, and to facilitate our study of subdivision of cones, we now introduce parallelotopes.

Definition 3B.2. If $\underline{a}_1, \dots, \underline{a}_n$ are linearly independent points in R^d , then the parallelotope of $\underline{a}_1, \dots, \underline{a}_n$ is defined as

$$P = P(\underline{a}_1, \dots, \underline{a}_n) = \{\lambda_1 \underline{a}_1 + \dots + \lambda_n \underline{a}_n \mid 0 \leq \lambda_i < 1 \forall i\}.$$

Note that P is a half-open parallelotope which does not include the vertices $\underline{a}_1, \dots, \underline{a}_n$. The following result is easy to see.

Proposition 3B.2. Let N, M be n -dimensional lattices with M a sublattice of N , let $\underline{a}_1, \dots, \underline{a}_n$ be a basis for M , and let P be the parallelotope of $\underline{a}_1, \dots, \underline{a}_n$. Then

- (i) $N \cap P$ is a complete set of representatives for the cosets of M in N .
- (ii) $N = M \Leftrightarrow (N:M) = 1 \Leftrightarrow N \cap P = \emptyset$.

Base points of, and notation for L-cones. If C is an n -dimensional L-cone in R^d with r edges, then C is determined by the r primitive points $\underline{a}_1, \dots, \underline{a}_r$ of L on the edges of C .

Definition 3B.3. The points a_1, \dots, a_r are called a base for C , and to indicate this, we will use the notation

$$C = c(a_1, \dots, a_r).$$

More generally, the primitive points on the edges of a collection S of L -cones will be called the base points of S , and $B(S)$ will denote the set of all these points.

Index of a simplicial L -cone.

Definition 3B.4. If $C = c(a_1, \dots, a_n)$ is a simplicial L -cone, the index of C (in L) is the index $(L':M)$, of the sublattice $M = \ell(a_1, \dots, a_n)$ in $L' = L \cap \text{lin } C$, and is denoted by $m(C)$.

Thus, C is a basic L -cone (as defined in §3A.1) if and only if its index is 1. Notice that $\dim L' = \dim C = n$, and that the index of C in L is equal to its index in L' . We do not suppose that $L' = L$, that is $\dim L = \dim C$, as we wish to be able to consider a complex of L -cones C such that the $\text{lin } C$ do not all coincide.

Definition 3B.5. If C is a simplicial L -cone, then the parallelotope of C is the parallelotope $P(a_1, \dots, a_n)$ of its base. Denote it by $P(C)$.

From Proposition 3B.2 we now have the following result in terms of $P(C)$ (which is equivalent to Lemma 3A.3 together with its converse).

Lemma 3B.3. If C is a simplicial L -cone, then C is basic if and only if $L \cap P(C) = \emptyset$ (that is, if and only if the only points of L in the closure $\bar{P}(C)$ of $P(C)$ are the 2^n vertices).

We remark that if C is a simplicial L -cone, and \tilde{x} is indecomposable in $L \cap C$, then \tilde{x} is either a base point of C , or belongs to $P(C)$.

Outer index of a general L -cone. This is a notion which will be used in expressing bounds in Theorem 3C.2. (The notion of inner index will be introduced in §3D.)

Definition 3B.6. If $C = c(\tilde{a}_1, \dots, \tilde{a}_r)$ is an n -dimensional L -cone, the (outer) index of C (in L) is

$$m(C) = v(C)/v_0(C),$$

where $v(C)$ is the (n -dimensional) volume of

$$(1) \quad S(C) = \text{conv}(\tilde{0}, \tilde{a}_1, \dots, \tilde{a}_r),$$

and $v_0(C)$ is the volume of the simplex generated by $\tilde{0}$ and any basis of the lattice $L' = L \cap \text{lin } C$.

Clearly, if C is simplicial, then its index (as defined in Definition 3B.4) is equal to its outer index. (The index is also equal to the inner index defined later in §3D.)

Lemma 3B.4.

Let S be a complex of n -dimensional L -cones in \mathbb{R}^d .

Then

- (i) By repeated barycentric subdivision, we may obtain a simplicial subdivision S' of S (see Note 3A.1.) such that the base points of S' are precisely those of S .
- (ii) Any subdivision S_1 of S (whether simplicial or not) whose base points all lie in the union of the $S(C)$ given by (1) for $C \in S$ satisfies

$$(2) \quad \sum_{D \in S_1} m(D) \leq \sum_{C \in S} m(C)$$

Proof. (i) is just Lemma 3A.2, expressed in terms of base points. (ii) is clear from volume considerations for C in S .

It can be shown that by making a subdivision S_1 of S whose cones correspond to the facets of the $S(C)$ not through Q , and then using barycentric subdivision, one may obtain equality in (2). Another 'facet construction' is considered later in §3D, where we will minimise, rather than maximise the left hand side of (2) (thus obtaining what will be called a minimal special subdivision of S).

Special cones.

We now introduce the notion of special cones in terms of deleted simplexes. Special cones will result from the construction in §3D, and are needed for Chapters 4 and 5.

Definition 3B.7. If $S = \text{conv}(\underline{a}_0, \dots, \underline{a}_n)$ is the simplex with the $n+1$ affinely independent points $\underline{a}_0, \dots, \underline{a}_n$ as vertices, then the deleted simplex is

$$S^*(\underline{a}_0, \dots, \underline{a}_n) = S \sim \{\underline{a}_0, \dots, \underline{a}_n\}$$

Definition 3B.8. If $C = c(\underline{a}_1, \dots, \underline{a}_n)$ is an L-cone, then the deleted simplex $S^*(C)$ is

$$S^*(C) = S^*(Q, \underline{a}_1, \dots, \underline{a}_n) = \{\sum \lambda_i \underline{a}_i \in P(C) \sim \{Q\} \mid \sum \lambda_i \leq 1\}.$$

Definition 3B.9. If $C = c(\underline{a}_1, \dots, \underline{a}_n)$ is a simplicial L-cone, then C is special if

$$L \cap S^*(C) = \emptyset,$$

that is, if

$$\left. \begin{array}{l} \underline{a} = \lambda_1 \underline{a}_1 + \dots + \lambda_n \underline{a}_n \in L \cap C, \\ \lambda_1 + \dots + \lambda_n \leq 1 \end{array} \right\} \Rightarrow \begin{array}{l} \underline{a} = Q, \text{ or} \\ \underline{a} = \underline{a}_i \text{ for} \\ \text{some } i. \end{array}$$

Lemma 3B.5. A 2-dimensional L-cone C is special if and only if it is basic.

Proof. Let $C = c(\underline{a}_1, \underline{a}_2)$ be a 2-dimensional simplicial cone. By Lemma 3B.3, C is basic if and only if $\dot{L} \cap P(C) = \emptyset$. Hence if C is basic, it is obviously special. Conversely, suppose C is special. To prove that C is basic, suppose to the contrary, that $\dot{L} \cap P(C) \neq \emptyset$, so that some

$$\underline{a} = \lambda_1 \underline{a}_1 + \lambda_2 \underline{a}_2 \in \dot{L} \cap P(C).$$

Then

$$\underline{a}' = (1-\lambda_1)\underline{a}_1 + (1-\lambda_2)\underline{a}_2 \in \dot{L}.$$

If $\lambda_1 + \lambda_2 \leq 1$, then $\underline{a} \in L \cap S^*(C)$, and if $\lambda_1 + \lambda_2 \geq 1$, then $\underline{a}' \in L \cap S^*(C)$. This contradicts the assumption that C is special.

Lemma 3B.5 is Lemma 6, of Chapter III of Cassels [3], and is also (essentially) Theorem 34 of Hardy and Wright [6]. The result does not hold for n -dimensional cones with $n > 2$. In Chapter 5 on 3-dimensional L -cones we will obtain and use a characterisation of special cones.

It is clear that if C is a special L -cone, then any face of C (of dimension at least one) is also special.

§3C. Basic subdivisions of a complex of L-cones.

In this section, L will again be a fixed lattice in R^d of dimension $\leq d$. We will solve the subdivision problem, and provide the relevant bounds. Then we will give applications to points of a sublattice of Z^n which lie in the non-negative orthant, and (Corollary 3C.3) solve Schinzel's problem, in its original form.

The main theorem of this section, Theorem 3C.2 on subdivision, is based on the following key lemma on barycentric subdivision of a complex of simplicial L-cones.

Lemma 3C.1. Let S be a complex of n -dimensional simplicial L-cones in R^d . Let C be any cone in S of index $m(C) > 1$. Then the following hold.

- (i) There is a point \tilde{a} in $\dot{L} \cap P(C)$.
- (ii) If D is any cone in S containing the point \tilde{a} (of $\dot{L} \cap P(C)$), then $\tilde{a} \in \dot{L} \cap P(D)$, and the index $m(D) > 1$.
- (iii) In barycentric subdivision of S with division point \tilde{a} , the (at most n) cones D_i replacing D all have index less than $m(D)$.

Proof (i) This is immediate from Lemma 3B.3.

- (ii) Since S is a complex, \tilde{a} belongs to the common face $F = C \cap D$. But F is a simplicial L-cone, and its base is a subset of the base of C , and is also a subset of the base of D . Thus $\tilde{a} \in P(C) \cap F = P(F) = P(D) \cap F \subseteq P(D)$. This shows that $\tilde{a} \in \dot{L} \cap P(D)$, and so by Lemma 3B.3 again, $m(D) > 1$.

(iii) The lattice point \underline{a} is given by

$$0 \neq \underline{a} = \lambda_1 \underline{a}_1 + \dots + \lambda_n \underline{a}_n,$$

where $0 \leq \lambda_i < 1$ for all i , and $\underline{a}_1, \dots, \underline{a}_n$ is the base of D . In barycentric subdivision, each cone D is replaced by the cones

$$D_i = c(\underline{a}_1, \dots, \underline{a}_{i-1}, \underline{a}, \underline{a}_{i+1}, \dots, \underline{a}_n), \text{ where } \lambda_i > 0.$$

The cone D_i has index $\lambda_i m(D)$ by definition of the index of a simplicial L-cone (Definition 3B.4.) and by Proposition 3B.1. Thus D is replaced by at most n cones D_i , all of index less than $m(D)$. This completes the proof of Lemma 3C.1.

If S is a collection of L-cones, let

$$N(S) = \#\{C \mid C \in S\},$$

$$E(S) = \#\{\underline{b} \mid \underline{b} \text{ is a base point of } S\},$$

$$H(S) = \max\{H(\underline{b}) \mid \underline{b} \text{ is a base point of } S\},$$

where $H(\underline{b})$ is the height of \underline{b} (see definition in §3A.3), and let $m(C)$ denote the outer index of an L-cone C (Definition 3B.6). Notice that these numbers correspond respectively to (ii), (i), (iii) in §3A.4.

Theorem 3C.2. (Subdivision Theorem.)

Let S be a complex of n -dimensional L -cones in R^d , where $n \geq 2$. Let S_1 be (as in Lemma 3B.4) a simplicial subdivision of S such that the base points of S_1 all lie in the union of the $S(C)$ for C in S . Then we obtain a basic subdivision T of S_1 (and so of S) by the following process:

Perform a sequence of barycentric subdivisions, where at each step the division point may be chosen to be any point in $\dot{I} \cap P(C)$, where C is any (simplicial) L -cone of the subdivision reached with index $m(C)$ greater than 1, and stop when all cones of the subdivision reached are basic.

Then the inequalities below apply, where $f(a,b)$ denotes the generalised Fibonacci sequence defined by

$$\begin{cases} f(a,b) = 1 & \text{for } b = -a + 1, -a + 2, \dots, 0, \\ f(a,b) = f(a,b-a) + \dots + f(a,b-1) & \text{for } b \geq 1, \end{cases}$$

and $N(S)$, $E(S)$, $H(S)$ are defined as above.

For the simplest case where S consists of just one simplicial L -cone C of index m , we have

$$(i) \quad N(T) \leq n^{m-1}$$

$$(ii) \quad E(T) \leq n + \frac{n^{m-1}-1}{n-1} = E(C) + \frac{n^{m-1}-1}{n-1}$$

$$(iii) \quad H(T) \leq H(C) f(n, m-1).$$

In the case of arbitrary S , we have, on letting

$$N' = \sum_{C \in S} n^{m(C)-1},$$

$$(i)' \quad N(T) \leq N'$$

$$(ii)' \quad E(T) \leq E(S) + \frac{1}{n-1} (N' - N(S))$$

$$(iii)' \quad H(T) \leq \max_{C \in S} (H(C) f(n, m(C)-1))$$

Proof. (We note that there is no subdivision problem for $n = 1$, as all 1-dimensional L-cones are basic.) The key Lemma 3C.1 shows that the processes described above all terminate with basic L-cones, and ~~leads to~~ the bounds in (i), (ii), (iii). Parts (i)' to (iii)' follow easily. The proof of (i)' uses Lemma 3B.4(ii) and uses the inequality

$$n^{m_1-1} + \dots + n^{m_s-1} \leq n^{m_1+\dots+m_s-1} \quad (m_i \geq 1 \forall i, \text{ and } n \geq 2).$$

The proof of (ii)' uses the fact that $N(S_1) \geq N(S)$.

Corollary 3C.3. (Solution of the special case of Schinzel's problem of §3A.2.)

Let L be an n -dimensional sublattice of \mathbb{Z}^n such that the index $(\mathbb{Z}^n : L) = r$. Then the non-negative orthant

$$E^+ = \{\underline{x} = (x_i) \in \mathbb{R}^n \mid x_i \geq 0 \forall i\}$$

has a subdivision T into basic L-cones such that

$$\begin{aligned} N(T) &\leq nr^{n-1-1}, \\ E(T) &\leq n + \frac{nr^{n-1-1}}{n-1}, \\ H(T) &\leq rf_0(n, r^{n-1}-1), \end{aligned}$$

where $f_0(a, b)$ is the generalised Fibonacci sequence $f(a, b)$ modified (diminished) by taking the initial values to be $0, \dots, 0, 1$.

Proof. By Cramer's rule, the points re_i , where e_i is the i th natural unit vector, belong to L . Thus, E^+ is an L-cone. Its primitive base points are r_1e_1, \dots, r_ne_n , where $r_i \geq 1$, and $r_i \mid r$ for $i = 1, \dots, n$. Also,

$$M = \ell(r_1e_1, \dots, r_ne_n) \subseteq L \subseteq \mathbb{Z}^n.$$

The index of E^+ in L is defined to be the index $(L:M)$.

But

$$(Z^n:M) = (Z^n:L) (L:M), \text{ or}$$

$$r_1 \dots r_n = r(L:M).$$

Hence $(L:M)$, equal to m say, divides r^{n-1} . The result now follows by applying the theorem, with f_0 replacing f in part (iii) to fit in with the special circumstances here, when all but one of the coordinates of each base point of S are zero. (The modified part (iii) is seen in the same way as the original.)

Corollary 3C.4. (Generalisation of Schinzel's problem to non-full lattices in R^d .) Let L be a sublattice of Z^d of dimension possibly less than d . Let E^+ be the non-negative orthant in R^d , let $L^+ = L \cap E^+$, and let $\dim L^+ = n$. Then there is a finite subset B of L^+ such that every point x in L^+ is expressible as a non-negative integral linear combination of an n -element subset of B .

Proof. Let $C = (\text{lin } L^+) \cap E^+$ (which can easily be shown to equal $(\text{lin } L) \cap E^+$). The result will follow once we have shown that

$$(i) \quad L^+ = L \cap C$$

$$(ii) \quad C \text{ is an } n\text{-dimensional } L\text{-cone.}$$

$$\begin{aligned} \text{Proof of (i) : } L^+ &= L^+ \cap \text{lin } L^+ = (L \cap E^+) \cap \text{lin } L^+ \\ &= L \cap ((\text{lin } L^+) \cap E^+) \\ &= L \cap C. \end{aligned}$$

Proof of (ii) : The fact that the polyhedral cone C is pointed with apex 0 follows from the fact that E^+ is pointed with apex 0 . Since $L^+ \subseteq C \subseteq \text{lin } L^+$, C is n -dimensional. There remains to prove that each edge of

C contains a point of L . Now C is defined by inequalities with rational coefficients. Hence each edge of C contains a point a of $Z^d \cap \text{lin } L^+$, and to prove that it contains a point ra of L ($r \in \mathbb{Z}$), it is sufficient to prove that both $Z^d \cap \text{lin } L^+$ and the sublattice $L \cap \text{lin } L^+$ are n -dimensional. But $\text{lin } L^+$ has a basis of n points in L^+ , and $L^+ \subseteq L$, so $\dim(L \cap \text{lin } L^+) = n$, and $Z^d \cap \text{lin } L^+ \subseteq \text{lin } L^+$, so $\dim(Z^d \cap \text{lin } L^+) \leq n$, and the result follows.

Note (added August, 1979).

The above corollary can be used to obtain a refinement of Theorem 2 of the paper

A. Schinzel, Reducibility of lacunary polynomials. I,
Acta Arith. 16 (1969), 123-159.

§3D. Minimal special subdivisions, and inner index of a complex of L-cones.

The bounds obtained in Theorem 3C.2 are quite crude ones. We shall find sharp bounds in Chapters 4, 5 for 2, 3-dimensional cones. The method will use the results in the present section, which concern a choice of the S_1 in Theorem 3C.2, and involve the facets of the convex hull of the non-zero lattice points in an L-cone.

It is plausible that to obtain an "economical" subdivision, that is, one with low bounds, one should start off the subdivision of S by selecting new base points from within the sets $S(C)$ (see §3B(1)) for C in S , while minimising the sum of the outer indices (see Definition 3B.6) of the resulting cones. New base points outside any $S(C)$ are then used only if any non-basic cones result. In fact, this idea does work in 2 and 3 dimensions, as we will see in the applications of this section to Chapters 4, 5. The lemma and theorem below are fairly obvious results in 2, 3 dimensions.

Lemma 3D.1. Let C be an n -dimensional L-cone in \mathbb{R}^d ($n \geq 2$), let $I(C)$ denote (as usual) the set of points indecomposable in $\dot{L} \cap C$, and let

$$V(C) = U\{F \mid F \text{ is a facet of } \text{conv}(\dot{L} \cap C), \text{ and } 0 \notin \text{aff } F\}.$$

Then

- (i) $\text{conv}(\dot{L} \cap C)$ is an n -dimensional convex polyhedral set and does not contain 0 .
- (ii) Each facet F of $\text{conv}(\dot{L} \cap C)$ such that $0 \notin \text{aff } F$ is an $(n-1)$ -dimensional polytope whose vertices are points of L .

(iii) The set $V(C)$ consists precisely of the points of $\text{conv}(\dot{L} \cap C)$ which are visible from 0 (that is, points \underline{x} in $\text{conv}(\dot{L} \cap C)$ such that the line segment $0\underline{x}$ contains no point of $\text{conv}(\dot{L} \cap C)$ other than \underline{x}). Every ray in C from 0 passes through exactly one point of $V(C)$.

(iv) Let F be a facet of $\text{conv}(\dot{L} \cap C)$ such that $0 \notin \text{aff } F$. Let $\underline{n} = \underline{n}(F)$ be the unique vector such that

$$\underline{n} \cdot \underline{x} = 1 \quad \forall \underline{x} \in F, \text{ and } \underline{n} \cdot \underline{x} \geq 1 \quad \forall \underline{x} \in \text{conv}(\dot{L} \cap C).$$

Then

$$\{\underline{x} \mid \underline{x} \in \dot{L} \cap C, \underline{n} \cdot \underline{x} < 2\} \subseteq I(C).$$

(v) $L \cap V(C) \subseteq I(C)$.

Proof. Parts (i) - (iii) are easily seen.

(iv) If $\underline{x} \in \dot{L} \cap C$, but $\underline{x} \notin I(C)$, then $\underline{x} = \underline{x}_1 + \underline{x}_2$, where $\underline{x}_1, \underline{x}_2 \in \dot{L} \cap C$, so $\underline{n} \cdot \underline{x} = \underline{n} \cdot \underline{x}_1 + \underline{n} \cdot \underline{x}_2 \geq 2$ by (iv).

(v) This follows from (iv), since if $\underline{x} \in L \cap V(C)$, then $\underline{x} \in L$, and $\underline{x} \in F$ for some facet F of $\text{conv}(\dot{L} \cap C)$ with $0 \notin \text{aff } F$.

In the following theorem, $m(D)$ will again denote the outer index of an L -cone D (as defined in Definition 3B.6).

Theorem 3D.2. Let S be a complex of n -dimensional L -cones in R^d ($n \geq 2$). For C in S , let $V(C)$ be as defined in Lemma 3D.1. Let

$$S_0 = \{\text{cone}_0(F) \mid 0 \notin \text{aff } F, \text{ and } F \text{ is a facet of } \text{conv}(\dot{L} \cap C) \text{ for some } C \in S\},$$

and let

$$T = \bigcup_{C \in S} (L \cap V(C)) = L \cap \bigcup_{C \in S} V(C).$$

Then

- (i) S_0 is a subdivision of S .
- (ii) There is a special subdivision S_1 of S_0 (that is, a subdivision into n -dimensional special L -cones) such that T is the set of base points of S_1 . (Special L -cones were defined in Definition 3B.9.)
- (iii) For any S_1 as in (ii), and any subdivision V of S ,

$$\sum_{D \in S_1} m(D) = \sum_{D \in S_0} m(D) \leq \sum_{D \in V} m(D) \leq \sum_{D \in S} m(D)$$

Proof. (i) It is easy to see that the $\text{conv}(\dot{L} \cap C)$ form a complex, and hence the F do too. Lemma 3D.1(iii) then implies that any ray of C from 0 passes through exactly one point of $\bigcup_{C \in S} V(C)$. It follows that S_0 is a complex. The other details are easily checked. (Use, for example, Lemma 3D(ii).)

(ii) A subdivision S_1 of the type required is obtained by simply applying repeated barycentric subdivision to S_0 with points of T as division points.

(iii) This follows easily by considering volumes, or indices, for each C in S , using the visibility result in part (iii) of the lemma for the first inequality.

We note that to economically subdivide a given complex, we could try forming S_1 as in the theorem, and then subdividing the resulting special cones into basic cones. This is what we do in Chapters 4 and 5, although in the 2-dimensional case (that is, Chapter 4), the second step is trivial.

Definition 3D.1. A subdivision of the type S_1 in Theorem 3D.2 of a complex S of n -dimensional L -cones will be called a minimal special subdivision of S . The integer

$$\sum_{D \in S_1} m(D) = \sum_{D \in S_0} m(D)$$

will be called the inner index of S .

Notation. We denote the inner index of a complex S of n -dimensional L -cones by

$$m_i(S) .$$

Thus, for simplicial C ,

$$m_i(C) = m_i(\{C\}) = m(C) ,$$

that is, the inner index coincides with the outer index (and the ordinary index). It is also clear that

$$m_i(S) = \sum_{D \in S} m_i(D) .$$

Remark. We have already remarked that the bounds obtained in Theorem 3C.2 are crude ones. However, to illustrate the concept of the inner index, we note that (as is easily seen) the result (i)' therein may be sharpened to

$$N(T) \leq n^{m_i(C)-1}$$

in the case when $S = \{C\}$ consists of a single L -cone.

§3E. Results on Indecomposable points.

In this section we gather together for future reference some results on the set $I = I(C)$ of indecomposable points of an L-cone C .

As already mentioned, the points of this set either belong to $P(C)$ or are base points of C , and they are included amongst the base points of any subdivision (or even decomposition) of C into basic L-cones.

The proof of the next result will use a result from Chapter 2 on residues (Lemma 2.2(vi)).

Proposition 3E.1. Let $C = c(\underline{a}_1, \dots, \underline{a}_n)$, $n \geq 2$, be a simplicial L-cone in R^d , and let

$$(1) \quad \underline{a} = \frac{1}{m} (\alpha_1 \underline{a}_1 + \dots + \alpha_n \underline{a}_n)$$

be in $L \cap C$, where the α_i are non-negative integers.

Suppose that

$$(2) \quad \alpha_1 + \dots + \alpha_n > (n-1)m - (n-2).$$

Then \underline{a} is decomposable in $L \cap C$.

Proof. The result is trivial for $m = 1$. So suppose $m \geq 2$. Since the right hand side of (2) is greater than or equal to m , $\underline{a} \neq \underline{a}_i$ for any i . Hence, if some $\alpha_i \geq m$, then $\underline{a} = \underline{a}_i + (\underline{a} - \underline{a}_i)$ is a non-trivial decomposition of \underline{a} , and so \underline{a} is decomposable. Also, if some $\alpha_i = 0$, then since the right hand side of (2) is equal to $(n-1)(m-1) + 1$, some $\alpha_j \geq m$, so by the previous case, \underline{a} is again decomposable. Thus, there remains only the case when $1 \leq \alpha_i \leq m-1$ for $i = 1, \dots, n$. By Lemma 2.2(vi), there is an integer r such that $2 \leq r \leq m-1$, and $\overline{r\alpha_i} < \alpha_i$ for $i = 1, \dots, n$, where \overline{x} denotes the least non-negative residue of $x \bmod m$. Let

$$\underline{b} = \frac{1}{m} (\overline{r\alpha_1} \underline{a}_1 + \dots + \overline{r\alpha_n} \underline{a}_n), \text{ and } \underline{c} = \underline{a} - \underline{b}. \text{ Then}$$

$b, c \in L \cap C$, and $c \neq 0$. Now consider the special case when the greatest common divisor of $\alpha_1, \dots, \alpha_n, m$ is 1. Then $b \neq 0$, so $a = b + c$ is a non-trivial decomposition, showing that a is decomposable. In the general case when the greatest common divisor is $d > 1$, (2) gives

$$\frac{\alpha_1}{d} + \dots + \frac{\alpha_n}{d} > (n-1)\frac{m}{d} - \frac{n-2}{d} \geq (n-1)\frac{m}{d} - (n-2),$$

and (1) gives

$$a = \frac{1}{m/d} \left(\frac{\alpha_1}{d} a_1 + \dots + \frac{\alpha_n}{d} a_n \right),$$

and indecomposability of a follows from the result in the special case applied to a with m/d in place of m . This completes the proof of the proposition.

We know that every point a of $L \cap C$ is of the form (1), where m is the index of C . Of particular interest in Chapters 4, 5 are the cases $n = 2, 3$, which yield the following corollaries of Proposition 3E.1.

Corollary 3E.2. Let $C = c(a_1, a_2)$ be a 2-dimensional L -cone, and let $a = \lambda_1 a_1 + \lambda_2 a_2 \in \dot{L} \cap C$. If $\lambda_1 + \lambda_2 > 1$, then a is decomposable.

Corollary 3E.3. Let $C = c(a_1, a_2, a_3)$ be a 3-dimensional L -cone, and let $a = \lambda_1 a_1 + \lambda_2 a_2 + \lambda_3 a_3 \in \dot{L} \cap C$. If $\lambda_1 + \lambda_2 + \lambda_3 \geq 2$, then a is decomposable.

Finally, we give a result on indecomposable points of complexes, which will be used in §5D.

Lemma 3E.4. (i) Let $a \in \dot{L} \cap F$, where F is a face of the L -cone C . Then a is indecomposable in C if and only if a is indecomposable in F .

(ii) Let S be a complex of L -cones in R^d , let $C, D \in S$, and let $a \in \dot{L} \cap C \cap D$. Then a is

indecomposable in C if and only if a is
indecomposable in D .

Proof. (i) Obvious.

(ii) This follows from (i).

Thus, we make the following definition.

Definition 3E.1. If S is a complex of L -cones, then an indecomposable point of S is an indecomposable point of a member of S . The notation $I(S)$ will be used for the set of indecomposable points of S .

Note 3E.5. If S is a basic complex, that is, a complex of basic cones, and $B(S)$ is the set of base points of S , then $B(S) = I(S)$. For general S , $B(S) \subseteq I(S)$. If T is a subdivision of S , then $I(S) \subseteq I(T)$.

CHAPTER 4

SUBDIVISION OF 2-DIMENSIONAL LATTICE CONESIntroduction

In this chapter I will consider fairly completely the subdivision problem for a 2-dimensional L-cone C in R^d . A subdivision of C into basic (2-dimensional) L-cones will be called a basic subdivision of C .

After giving some preliminary lemmas in §4A, I will show, in §4B, that C has a unique basic subdivision T_0 whose base points are indecomposable in $\dot{L} \cap C$, and hence I will derive bounds for Schinzel's problem in this case. In §4C I will consider other basic subdivisions of C , and show that any basic subdivision is obtained from the "minimal" basic subdivision T_0 by constructing a sequence of "mediants". In §4D I will relate the notion of subdivision of 2-dimensional L-cones to Farey sequences, convergents and best approximations; and in §4E I will apply the ideas of this chapter to a lemma of Schinzel.

§4A. Preliminary.

In this section I will give some preliminary lemmas, and I will discuss the relation between decompositions, dissections and subdivisions into arbitrary cones, basic cones and basic cones with indecomposable base points. Hence it will appear that there is no real loss of generality in investigating basic subdivisions of a 2-dimensional L-cone, rather than basic decompositions.

Any 2-dimensional L-cone is clearly simplicial, that is of the form

$$C = c(\underline{a}_1, \underline{a}_2),$$

where $\underline{a}_1, \underline{a}_2$ are the base points of C (see Definition 3B.3).

If C_1, C_2 are two 2-dimensional L-cones lying in the same plane, then clearly a necessary and sufficient condition that C_1, C_2 have no relative interior point in common is that $C_1 \cap C_2$ is a face of C_1, C_2 (which may be an edge, or \emptyset). Hence the concepts of dissection and subdivision of a 2-dimensional L-cone (see §3A.3) are identical.

We next consider the relation between dissections and decompositions, assuming that our subcones are basic. For this we need the following lemma.

Lemma 4A.1. Let $C_1 = c(\underline{a}_1, \underline{a}_2), C_2 = c(\underline{b}_1, \underline{b}_2)$ be basic 2-dimensional L-cones included in some 2-dimensional cone C . Then either one of C_1, C_2 is contained in the other, or their relative interiors are disjoint.

Proof. Suppose to the contrary that C_1, C_2 overlap (that is, their relative interiors are not disjoint), and that neither of C_1, C_2 is contained in the other. Then the four edges of C_1, C_2 are distinct, and, without loss of generality, we may suppose that they are ordered by angular displacement in the order $\underline{0a}_1, \underline{0b}_1, \underline{0a}_2, \underline{0b}_2$, where the last three rays are in the same one of the two open half-planes of $\text{lin } C$ bounded by the line $\underline{0a}_1$. Let φ be the linear transformation from $\text{lin } C$ to R^2 for

which

$$\varphi(\underline{a}_1) = \underline{e}_1, \quad \varphi(\underline{a}_2) = \underline{e}_2,$$

where $\underline{e}_1, \underline{e}_2$ are the unit vectors, and suppose that

$$\varphi(\underline{b}_1) = \underline{c}_1, \quad \varphi(\underline{b}_2) = \underline{c}_2,$$

so that

(1) The points $\underline{c}_1, \underline{c}_2$ are in the first and second open quadrants respectively.

Since C_1, C_2 are basic, the base points of $\varphi(C_1), \varphi(C_2)$ generate the same lattice. Hence

$$(2) \quad \underline{c}_1, \underline{c}_2 \in \mathbb{Z}^2, \quad \det(\underline{c}_1, \underline{c}_2) = 1.$$

Now (1), (2) lead to a contradiction, for they imply that for positive integers e_1, f_1, e_2, f_2 ,

$$\underline{c}_1 = (e_1, f_1), \quad \underline{c}_2 = (-e_2, f_2),$$

$$\begin{vmatrix} e_1 & f_1 \\ -e_2 & f_2 \end{vmatrix} = 1.$$

But the value of the determinant is $e_1 f_1 + e_2 f_2$, which is at least 2. This contradiction completes the proof of the lemma.

Note. The condition that C_1, C_2 are both contained in some cone C is essential. For example, if $\underline{a}_1, \underline{a}_2$ were both in the relative interior of C_2 , and $C_3 = c(-\underline{a}_1, \underline{a}_2)$, then C_3, C_2 would not satisfy the condition, nor the conclusion of the lemma.

Corollary. Let C be a 2-dimensional L-cone, and let

$$C = C_1 \cup \dots \cup C_r,$$

where the C_i are 2-dimensional basic L-cones, and the C_i are all irredundant, that is,

$$C_i \not\subseteq \bigcup_{j \neq i} C_j \quad \text{for } i = 1, \dots, r.$$

Then the decomposition $\{C_1, \dots, C_r\}$ of C is a subdivision of C .

Proof. If

$$(3) \quad \text{relint } C_i \cap \text{relint } C_j \neq \emptyset,$$

then by the lemma, we may suppose that $C_i \subseteq C_j$, hence C_i is redundant. Thus (3) does not hold for any i, j . So $\{C_1, \dots, C_r\}$ is a dissection, and hence a subdivision of C , since C is 2-dimensional.

The following lemma, which is stated in a form convenient for use in Chapter 5, shows that the irredundancy condition in the above corollary is automatically satisfied if the base points of the C_i are all indecomposable.

Lemma 4A.2. Let C_1, C_2 be basic 1- or 2-dimensional L-cones, and let C be a 2-dimensional L-cone such that

$$C_1, C_2 \subseteq C,$$

$$B(C_1), B(C_2) \subseteq I(C)$$

(that is, the base points of C_1, C_2 are indecomposable in C). Then $C_1 \cap C_2$ is a face of C_1, C_2 .

Proof. The result is trivial unless at least one of C_1, C_2 is 2-dimensional, so suppose C_2 is 2-dimensional. The result will follow if we can show that no base point of C_1 can lie in $\text{relint } C_2$. Suppose, to the contrary, that

$$\underline{a} \in B(C_1) \cap \text{relint } C_2.$$

Then

$$\underline{a} \in I(C) \cap C_2 \subseteq I(C_2),$$

so $\underline{a} \in I(C_2)$. Since C_2 is basic, it follows that

$\underline{a} \in B(C_2)$, contradicting the fact that $\underline{a} \in \text{relint } C_2$.

The next lemma is an obvious result on arbitrary subdivisions of a 2-dimensional L-cone.

Lemma 4A.3. Let $C = c(\underline{a}_1, \underline{a}_2)$ be an L-cone, and let $\underline{b}_1, \dots, \underline{b}_r$ be r distinct primitive points in C , with $\underline{b}_1 = \underline{a}_1$, $\underline{b}_r = \underline{a}_2$, ordered according to increasing angular displacement from $\underline{0}\underline{a}_1$ to $\underline{0}\underline{a}_2$. Then there is a unique subdivision T of C (containing $r - 1$ cones) for which

$$B(T) = \{\underline{b}_1, \dots, \underline{b}_r\},$$

namely

$$T = \{D_1, \dots, D_{r-1}\},$$

where

$$D_i = c(\underline{b}_i, \underline{b}_{i+1}) \quad \text{for } i = 1, \dots, r-1.$$

We recall that by Lemma 3B.5, a 2-dimensional L-cone is basic if and only if it is special.

Suppose that C is a 2-dimensional L-cone of index $m(C) = m$. Since linear transformations do not affect the subdivision problem, we may suppose for the purpose of presenting examples (even though this does not help in the theory of subdivision), that

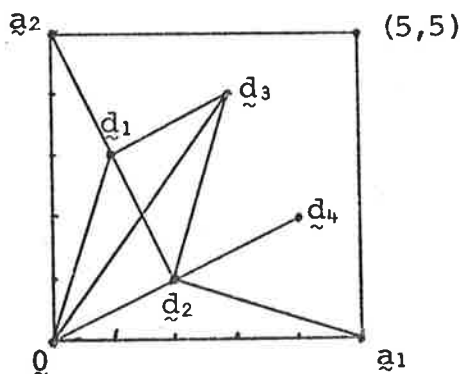
$$C = c(m\mathbf{e}_1, m\mathbf{e}_2), \quad m\mathbb{Z}^2 \subseteq L \subseteq \mathbb{Z}^2,$$

$$(L : m\mathbb{Z}^2) = (\mathbb{Z}^2 : L) = m.$$

We conclude this section by giving a numerical example of basic subdivisions of a 2-dimensional L-cone.

Example. Let $\underline{a}_1 = (5, 0) = 5\mathbf{e}_1$, $\underline{a}_2 = (0, 5) = 5\mathbf{e}_2$, $C = c(\underline{a}_1, \underline{a}_2)$, and let L be the lattice in \mathbb{Z}^2 generated by $5\mathbb{Z}^2$ and the point $(1, 3)$. Then the cone C is the

non-negative orthant E^+ in R^2 , and the index of C in L is $m(C) = m = 5$.



The paralleloptope $P(\underline{a}_1, \underline{a}_2)$ contains the following non-zero points of L :

$$\underline{d}_1 = (1, 3)$$

$$\underline{d}_2 = (2, 1)$$

$$\underline{d}_3 = (3, 4)$$

$$\underline{d}_4 = (4, 2).$$

Of these, $\underline{d}_4 = 2\underline{d}_2$ is not primitive, and $\underline{d}_3 = \underline{d}_1 + \underline{d}_2$ is decomposable in C . One basic subdivision T_0 of C is the subdivision (see Lemma 4A.3) for which

$$B(T_0) = \{\underline{a}_1, \underline{d}_2, \underline{d}_1, \underline{a}_2\} = I(C).$$

Another basic subdivision T_1 of C is given by

$$B(T_1) = \{\underline{a}_1, \underline{d}_2, \underline{d}_3, \underline{d}_1, \underline{a}_2\},$$

and T_1 is a subdivision of T_0 .

(In examples such as this, we can check whether a cone is basic by using determinants and indices, and we can use Cor 3E.2 to prove indecomposability.)

§4B. The minimal basic subdivision of a 2-dimensional L-cone.

Let C be a 2-dimensional L-cone. If T is a basic subdivision of C , then $B(T) = I(T)$, and $I(T) \supseteq I(C)$. In this section I will derive and apply the unique basic subdivision T_0 for which $I(T_0) = I(C)$, or equivalently, for which $B(T_0) = I(C)$. This unique T_0 will be called the minimal basic subdivision of C .

Theorem 4B.1. (Minimal basic subdivision.)

Let $C = c(a_1, a_2)$ be an L-cone. Then the minimal special subdivision T_0 of C (see Definition 3D.1) is unique, is basic, and is the unique subdivision satisfying

$$(1) \quad B(T_0) = I(T_0) = I(C) = L \cap V(C),$$

where $V(C)$ is the set of points of $\text{conv}(\dot{L} \cap C)$ which are visible from 0 .

Proof. By Theorem 3D.2(ii) and Lemmas 3D.1(v) and 3B.5 there is a minimal special, and hence basic, subdivision T_0 of C such that

$$B(T_0) = L \cap V(C) \subseteq I(C).$$

Hence

$$B(T_0) = I(T_0) \supseteq I(C).$$

Thus (1) is proved. The uniqueness of the subdivision satisfying (1) follows from Lemma 4A.3, and the uniqueness of the minimal special subdivision of C follows from this.

The above theorem, which does not depend on Cor 3E.2 leads easily to an alternative proof of that corollary.

A technique for calculating T_0 . We now illustrate a technique for determining $I(C)$ and hence T_0 . Given C , as in the Example of §4A, we will find a sequence of basic

cones

$$(2) \quad c(\underline{a}_2, \underline{a}), \quad c(\underline{a}, \underline{b}), \quad c(\underline{b}, \underline{a}_1)$$

whose base points belong to $I(C)$. The cones in (2) then form T_0 , and the base points $\underline{a}_2, \underline{a}, \underline{b}, \underline{a}_1$ comprise $I(C)$.

The point \underline{a} is uniquely determined by the conditions that $c(\underline{a}_2, \underline{a})$ be basic, and $\underline{a} \in I(C)$, for these imply that

$$\det(\underline{a}_2, \underline{a}) = \det((0, 5), \underline{a}) = -5,$$

$$\underline{a} = (1, 3) + k(0, 5), \quad k \in \mathbb{Z},$$

where k is least such that

$$(1, 3) + k(0, 5) \in C = E^+.$$

(The theorem implies that for this k , \underline{a} is in fact indecomposable.) Hence $k = 0$, $\underline{a} = (1, 3)$.

Similarly, \underline{b} is uniquely determined by the conditions that $c(\underline{a}, \underline{b})$ be basic, and $\underline{b} \in I(C)$. These imply that

$$\underline{b} = -\underline{a}_2 + h\underline{a},$$

where $h \in \mathbb{Z}$ is least such that $-\underline{a}_2 + h\underline{a} \in C$, and this gives $h = 2$, $\underline{b} = (2, 1)$.

Proceeding in this way we obtain, at the next step, the indecomposable point \underline{a}_1 .

From the technique just illustrated, it is easy to see that the minimal basic subdivision T_0 of C contains all the basic cones whose base points are indecomposable.

A consequence of Theorem 4B.1 is the following theorem on bounds, where $N(S)$, $E(S)$, $H(S)$ are as defined in §3C.

Theorem 4B.2. Let $C = c(\underline{a}_1, \underline{a}_2)$ be an L-cone of index m . Then there is a basic subdivision T of C

for which

$$(i) \quad N(T) \leq m$$

$$(ii) \quad E(T) \leq m+1$$

$$(iii) \quad H(T) \leq H(C).$$

Proof. Let T be the minimal basic subdivision T_0 of C , so that $B(T) = I(C)$. But

$$(3) \quad I(C) \subseteq (\dot{L} \cap P(C)) \cup \{a_1, a_2\},$$

and $\dot{L} \cap P(C)$ has $m-1$ elements. Thus

$$E(T) = \# B(T) = \# I(C) = m+1,$$

giving (ii). Hence, by Lemma 4A.3,

$$N(T) = E(T) - 1 = m,$$

that is, (i) holds. From (3) and Cor 3E.2,

$$I(C) \subseteq \text{conv}(Q, a_1, a_2),$$

whence (iii) follows.

The above bounds are best possible for cones C such that $\dot{L} \cap P(C)$ is contained in the line segment $a_1 a_2$.

§4C. Basic subdivisions of a 2-dimensional L-cone.

In this section I will first show the existence of basic subdivisions T of the 2-dimensional L-cone C for certain prescribed $B(T)$. I will then show that any basic subdivision of C can be obtained from the minimal basic subdivision T_0 by repeated barycentric subdivision, where the division points are "mediants".

Theorem 4C.1. Let S be any bounded convex subset of a 2-dimensional L-cone $C = c(a_1, a_2)$ such that

$$S \supseteq \text{conv}(Q, a_1, a_2).$$

Let A be the set of all primitive points of L which

lie in S . Then there is a basic subdivision T of C such that

$$B(T) = I(T) = A.$$

Proof. Since S is bounded, A is finite. Let T be the unique subdivision of C (see Lemma 4A.3) such that $B(T) = A$. We now show that T is special. Suppose to the contrary that $D \in T$, where the deleted simplex $S^*(D)$ contains a point of L . Then $S^*(D)$ contains a primitive point \underline{a} of L . Such a point \underline{a} cannot belong to A , since its angular displacement about $\underline{0}$ is between those of the base points of D . But S is convex, so $S^*(D) \subseteq S$, and $\underline{a} \in A$, a contradiction. Thus T is special, and hence (by Lemma 3B.5) basic, as required.

Any basic subdivision T of a 2-dimensional L -cone C is a subdivision of the minimal basic subdivision T_0 , because

$$B(T) \supseteq I(C) = B(T_0).$$

Hence it is trivial that T may be obtained from T_0 by repeated barycentric subdivision. For more detailed information on possible choices of division points, we need the following lemma which we will prove from Lemma 4A.1.

Lemma 4C.2. Let C be a 2-dimensional cone, and let $D = c(\underline{a}, \underline{b})$ and E be basic 2-dimensional L -cones contained in C . If $E \not\subseteq D$, then

$$\underline{a} + \underline{b} \notin \text{relint } E.$$

Proof. We will prove the contrapositive. Suppose that

$$(1) \quad \underline{a} + \underline{b} \in \text{relint } E.$$

Consider the basic cones

$$D_1 = c(\underline{a}, \underline{a+b}), \quad D_2 = c(\underline{a+b}, \underline{b}),$$

which are obtained from D by barycentric subdivision with primitive division point $\underline{a+b}$. It follows from (1), that for $i = 1, 2$,

$$\text{relint } E \cap \text{relint } D_i \neq \emptyset, \quad E \not\subseteq D_i.$$

Applying Lemma 4A.1 to the two basic cones E, D_i contained in C , we see that $E \supseteq D_i$, for $i = 1, 2$. Hence $E \supseteq D$. This completes the proof of the lemma.

Definition 4C.1. If T is a basic subdivision of the 2-dimensional L-cone C , and $D = c(\underline{a}, \underline{b}) \in T$, then the sum $\underline{a+b}$ will be called a mediant of T . If U is any basic subdivision of T obtained from T by repeated barycentric subdivision with mediant as division points, we will say that U is obtained from T by constructing mediants.

Notice that a mediant of T cannot be a base point of T . Hence the process of constructing mediant enlarges the set of base points.

Theorem 4C.3. (The mediant construction.) Let T be a basic subdivision of, and T_0 the minimal basic subdivision of the 2-dimensional L-cone C . Then T may be obtained from T_0 by constructing mediant.

Proof. We have that $B(T_0) \subseteq B(T)$, and by Lemma 4A.3, any set of primitive points in C determines a subdivision of C . Thus, it suffices to prove that if S is a basic subdivision of C such that $S \neq T$, and

$$(2) \quad B(S) \subseteq B(T),$$

then some mediant of S lies in T . Since $S \neq T$, there exists D such that

$$D = c(\underline{a}, \underline{b}) \in S, \quad D \notin T.$$

Now let E be any cone in T . By (2), $E \not\supset D$, so by Lemma 4C.2, $\underline{a} + \underline{b} \notin \text{relint } E$. Thus $\underline{a} + \underline{b} \in B(T)$, as required.

§4D. Farey sequences, convergents and best approximations.

Farey sequences, convergents and best approximations can be considered in terms of cones.

Let us represent each real number θ by the ray $\underline{Q}(\theta, 1)$ in \mathbb{R}^2 emanating from \underline{Q} and passing through the point $(\theta, 1)$, and let us represent the point at infinity by the positive x-axis $\underline{Q}(1, 0)$. Then the set of rays representing $\mathbb{R} \cup \{\infty\}$ sweep out the upper half-plane

$$(1) \quad H = \{(x, y) \mid x, y \in \mathbb{R}, y > 0 \text{ or } (y = 0, x \geq 0)\}$$

(where the negative x-axis, excluding \underline{Q} , is not included).

If θ is rational, say $\theta = r/s$, where r, s are relatively prime, and $s > 0$, then θ may also be represented by the unique primitive point (r, s) on the ray $\underline{Q}(\theta, 1)$ or $\underline{Q}(r, s)$, and ∞ may be represented by the point $(1, 0)$ on the positive x-axis.

Let

$$D = c(\underline{a}, \underline{b}),$$

where $\underline{a} = (a_1, a_2)$, $\underline{b} = (b_1, b_2)$ are primitive points of \mathbb{Z}^2 lying in H , and $\det(\underline{a}, \underline{b}) < 0$. Let $\underline{d} = (d_1, d_2) \in H$. Then

$$(2) \quad \frac{a_1}{a_2} \leq \frac{d_1}{d_2} \leq \frac{b_1}{b_2} \Leftrightarrow \underline{d} \in D.$$

Also,

$$(3) \quad \underline{d} \in \mathbb{Z}^2 \cap D \Leftrightarrow \underline{d} = \lambda \underline{a} + \mu \underline{b},$$

where

$$\lambda \geq 0, \quad \mu \geq 0, \quad \lambda, \mu \in \mathbb{Z}.$$

Farey sequences. Let G_n be the sequence

$$\dots, -\frac{1}{n}, \frac{0}{1}, \frac{1}{n}, \dots, \frac{1}{1}, \frac{n+1}{n}, \dots$$

of reduced fractions of denominator not exceeding n , arranged in usual real number order, and let F_n be the subsequence

$$\frac{0}{1}, \frac{1}{n}, \dots, \frac{1}{1}.$$

Either sequence may be referred to as the Farey sequence of order n . Consider also G'_n which is obtained by taking reduced fractions of height rather than denominator not exceeding n . (The height of a fraction in F_n is equal to its denominator.)

The standard construction of F_n begins with F_1 and proceeds by construction of mediants $(a_1+b_1)/(a_2+b_2)$ of consecutive fractions until there is no pair of consecutive fractions whose mediant has denominator not exceeding n . (It is not actually necessary to pass through F_2, F_3, \dots on the way to F_n .)

In terms of cones, the standard construction of F_n , considered as a set of primitive points in H , corresponds to the construction, by constructing mediants (see Definition 4C.1), of the basic subdivision T of the basic cone $c((0,1), (1,1))$ for which $B(T) = F_n$. The validity of the construction for F_n is most easily seen from this viewpoint, and is obvious from (2) and (3).

Similarly, G_n may be constructed by subdividing the basic cones $c((-1,0),(0,1))$ and $c((0,1),(1,0))$, and the same applies for G'_n .

Convergents and best approximations. We will consider the convergents of the real number θ from the cones viewpoint. Let A be the set of all basic Z^2 -cones which lie in H and contain $(\theta,1)$ in their relative interiors. There are two cases to consider.

(a) The case when θ is irrational.

We define the convergents of θ to be the base points of cones in A . We will order these convergents and show that, apart from an initial segment, they are the convergents (principal and intermediate) as defined in Lang [7].

Let $C, C' \in A$. From Lemma 4A.1, one of C, C' is contained in the other, so suppose that

$$C' \subseteq C = c(\underline{a}, \underline{b})$$

and let

$$E = c(\underline{a}, \underline{a}+\underline{b}), \quad F = c(\underline{a}+\underline{b}, \underline{b})$$

be the basic cones obtained from C by the mediant construction. Then from Lemma 4A.1 again it follows that $C' \subseteq E$ or F . Now this E or F belongs to A . Thus under the ordering \supseteq , the immediate successor of C is whichever of E, F contains $(\theta,1)$. Consequently, the immediate predecessor of C exists, and is whichever of $c(\underline{a}, \underline{b}-\underline{a}), c(\underline{a}-\underline{b}, \underline{b})$ lies in H . (Notice that exactly one of $\underline{a}-\underline{b}, \underline{b}-\underline{a}$ lies in H .) Thus the cones in A form a chain under \supseteq . From (3), it follows that there are only a finite number of basic cones of A lying between C and C' . Thus the cones of A may be indexed by the

integers, and form a descending sequence

$$(4) \quad C_{-2} \supseteq C_{-1} \supseteq C_0 \supseteq C_1 \supseteq C_2 \supseteq \dots$$

[From (3), we also see that the height of the base points of C_i tend to infinity as i tends to ∞ . Since the C_i are basic, this implies that

$$\bigcap_i C_i = Q(\theta, 1).$$

Also, if $t \leq [0]$, then $c((t, 1), (1, 0)) \in A$. Hence

$$\bigcup_i C_i = H. \quad]$$

We now see that the convergents of θ , defined as the base points of C in A , form a sequence

$$(5) \quad a_{-\infty} \ll \dots \ll a_{-1} \ll a_0 \ll a_1 \ll a_2 \ll \dots \ll a_i \ll \dots$$

where \ll is defined by

$$(6) \quad a \ll b \Leftrightarrow b - a \in H \sim \{0\},$$

$a_{-\infty} = (1, 0)$, and a_i is the mediant of C_{i-1} . Let us fix the indices by choosing

$$C_0 = c([0], 1), (1, 0).$$

Writing $a_0 = ([0], 1)$, $a_i = (\alpha_i, \beta_i)$ ($i > 0$),

we see that the α_i/β_i for $i \geq 0$ are precisely the convergents to θ , as defined in Lang [7].

Let D_1 denote the cone with apex Q generated by $(\theta, 1)$ and $(-1, 0)$ and let D_2 denote the cone with apex Q generated by $(\theta, 1)$ and $(1, 0)$. (These are not Z^2 -cones.) Define the linear function L by

$$(7) \quad L(x, y) = y\theta - x.$$

Then for \tilde{x} not lying on the ray $Q(\theta, 1)$,

$$(8) \quad L(\tilde{x}) \begin{cases} > 0 & \text{if } \tilde{x} \in D_1 \\ < 0 & \text{if } \tilde{x} \in D_2 \end{cases}.$$

We define the principal convergents to θ as the convergents a_i for which $L(a_i)$, $L(a_{i+1})$ have opposite signs,

and we define the intermediate convergents to θ as the convergents \underline{a}_i for which $L(\underline{a}_i), L(\underline{a}_{i+1})$ have the same sign. Now if \underline{a}_i is any convergent, then from the mediants method of construction of the sequence (5) of convergents, there is a unique j for which

$$(9) \quad c(\underline{a}_i, \underline{a}_j) \in A, \quad j < i,$$

and for this j , we have

$$(10) \quad L(\underline{a}_i) L(\underline{a}_j) < 0, \quad \underline{a}_i + \underline{a}_j = \underline{a}_{i+1}, \quad \text{and}$$

$$(11) \quad \underline{a}_k \in \text{relint } c(\underline{a}_i, \underline{a}_j) \Leftrightarrow k \geq i+1.$$

Also, with i, j as in (9), and for $h = i, j$, define E_h to be the cone with base \underline{a}_h and whichever of $(1,0), (-1,0)$ lies on the same side of $Q(\theta, 1)$ as \underline{a}_h . Then (11) is equivalent to

$$(12) \quad k \leq i \Rightarrow \underline{a}_k \in E_i \cup E_j,$$

and we also have

$$(13) \quad j < k \leq i \Rightarrow \underline{a}_k \in E_i.$$

From (10) the following is immediate.

Proposition 4D.1. Given the convergent \underline{a}_i , let j be as in (9). Then \underline{a}_i is a principal convergent if and only if

$$|L(\underline{a}_i)| < |L(\underline{a}_j)|.$$

Let θ be an irrational number, and as before let A be the set of all basic \mathbb{Z}^2 -cones which lie in H and contain $(\theta, 1)$. Define \underline{p} in $\dot{\mathbb{Z}}^2 \cap H$ to be a best approximation to θ if

$$(14) \quad \underline{p}' \in \dot{\mathbb{Z}}^2 \cap H, \quad \underline{p}' \ll \underline{p} \Rightarrow |L(\underline{p})| < |L(\underline{p}')|,$$

where \ll is as defined in (6) and L is as in (7). This agrees with the usual definition except for approximations with denominator 1, and except for the least element $(1,0)$

of $\dot{Z}^2 \cap H$ in the sense of \ll , which is a best approximation according to our definition.

Proposition 4D.2. The best approximations to θ are just the principal convergents to θ .

Proof. From the mediants construction of the sequence (5) of convergents it is clear that, apart from the ray $Q(\theta, 1)$, the whole of H is covered by (an infinite number of) non-overlapping basic Z^2 -cones $c(\underline{a}, \underline{b})$ where $\underline{a}, \underline{b}$ are convergents on the same side of $Q(\theta, 1)$, that is $L(\underline{a})L(\underline{b}) > 0$. Let \underline{d} be a point of $\dot{Z}^2 \cap H$, lying in $c(\underline{a}, \underline{b})$, and such that $\underline{d} \neq \underline{a}, \underline{b}$. By (3),

$$\underline{d} = \lambda \underline{a} + \mu \underline{b}$$

where

$$\lambda \geq 0, \quad \mu \geq 0, \quad \lambda + \mu \geq 2.$$

But

$$L(\underline{d}) = \lambda L(\underline{a}) + \mu L(\underline{b}),$$

hence

$$(15) \quad |L(\underline{d})| > \min(|L(\underline{a})|, |L(\underline{b})|).$$

Also,

$$(16) \quad \underline{a} \ll \underline{d}, \quad \underline{b} \ll \underline{d}.$$

It follows from (15), (16) that \underline{d} is not a best approximation. Thus any best approximation is a convergent. It also follows from (15), (16) that for the convergent $\underline{p} = \underline{a}_i$ to be a best approximation, it is sufficient that (14) be satisfied with the restriction that \underline{p}' be a convergent. That is, \underline{a}_i is a best approximation if and only if

$$(17) \quad k < i \Rightarrow |L(\underline{a}_i)| < |L(\underline{a}_k)|.$$

Let $j(< i)$ be as in (9). Then one of $\underline{a}_i, \underline{a}_j$ is in D_1 , and the other in D_2 . From the mediants construction

of the \underline{a}_k , and from (8), the \underline{a}_k in D_1 have positive $L(\underline{a}_k)$ which decrease with k , and the \underline{a}_k in D_2 have negative $L(\underline{a}_k)$ which increase with k . Hence, in each of D_1 and D_2 , and so in each of E_i and E_j (see (12)), $|L(\underline{a}_k)|$ decreases with k . Thus

$$(18) \quad k < i, \quad L(\underline{a}_k)L(\underline{a}_i) > 0 \Rightarrow |L(\underline{a}_i)| < |L(\underline{a}_k)|,$$

$$(19) \quad k < j, \quad L(\underline{a}_k)L(\underline{a}_j) > 0 \Rightarrow |L(\underline{a}_j)| < |L(\underline{a}_k)|.$$

From (13) it follows that $k < j$ may be replaced by $k < i$ in (19). From this and (18), we have

$$(20) \quad k < i \Rightarrow \min(|L(\underline{a}_i)|, |L(\underline{a}_j)|) < |L(\underline{a}_k)|.$$

If \underline{a}_i is a principal convergent, then by Proposition 4D.1, $|L(\underline{a}_i)| < |L(\underline{a}_j)|$, so (20) yields (17), and \underline{a}_i is a best approximation. On the other hand, if \underline{a}_i is an intermediate convergent, then (17) is false for $k=j$, and so \underline{a}_i is not a best approximation.

(b) The case when θ is rational; minimal basic subdivisions.

Let $\theta = r/s$, where (r,s) is a primitive point of \mathbb{Z}^2 lying in H . We define the convergents to r/s as for the irrational case, with the addition of the point (r,s) itself. The results are similar to those for the irrational case, but the chain (4) of cones of A terminates with a smallest cone $C_f = c(\underline{a}, \underline{b})$ such that $\underline{a} + \underline{b} = (r,s)$ and

$$\cap A = \cap_i C_i = C_f.$$

We will now state the relation between these convergents and minimal basic subdivisions. (The proof is not hard, and will be omitted.) For simplicity, let us suppose that $r/s > 0$, so that (r,s) lies in the first quadrant

$c((0,1), (1,0))$. Let

$$F_1 = c((0,1), (r,s)), \quad F_2 = c((r,s), (1,0)).$$

Then for $i = 1, 2$ the convergents of r/s which lie in F_i are the indecomposable points of F_i , and so determine the minimal basic subdivision of F_i ,

§4E. Note on a lemma of Schinzel.

In this note, we will use Theorem 4B.2(iii) to simplify the proof of, and at the same time improve a constant in Lemma 5 of Schinzel [11], a copy of which is given in the Appendix. We now give the modification to the lemma, and then the modification to the proof.

Modification to the statement of the lemma:

In the bounds for v_i, μ_i in (40), and for $h(T)$ in (46), make the replacement

$$(*) \quad "(2d^2)!^2" \text{ by } "(2d^2)!"$$

Modification to the proof:

Delete the passage beginning with "Let ξ_1 be" in the second sentence of the second paragraph on page 16 and ending with " $v \geq 0$ " at the end of the first sentence of the second paragraph on page 17. Apply the replacement (*) to the formula in the second last line of the proof. For the deleted section substitute the following:

"Let ξ be the least positive integer such that $(\xi, 0) \in M$, and let η be the least positive integer such that $(0, \eta) \in M$. Then

$$(47) \quad 0 < \xi, \eta < (2d^2)!$$

Let N be the sublattice (or submodule) of M with basis $(\xi, 0), (0, \eta)$. By Theorem 4B.2(iii) of L. Low's thesis, the

M -cone with base points $(\xi, 0)$, $(0, \eta)$, that is, the non-negative orthant in R^2 , has a subdivision into basis M -cones whose primitive basis points have height not exceeding $(2d^2)!$. Since the point (n, m) has components which are non-negative integers, it must belong to one of the basic M -cones, say one with base points

$$(v_1, \mu_1), (v_2, \mu_2).$$

By interchanging the order of these points if necessary, we may arrange that the matrix

$$M = \begin{bmatrix} v_1 & \mu_1 \\ v_2 & \mu_2 \end{bmatrix}$$

satisfies (41). Also, (40) and (42) hold."

Comparison of the two methods.

We took a sublattice of M with diagonal basis $(\xi, 0)$, $(0, \eta)$ and subdivided without increasing the bounds on the height of base points. Schinzel took a triangular basis (ξ_1, η_1) , $(0, -\eta_2)$ for M . What he did is equivalent to forming subdivisions of the basic M -cone $c((\xi_1, \eta_1), (0, -\eta_2))$ by the mediants construction to construct the ray $Q(1, 0)$. This is equivalent to the continued fraction process for obtaining η_1/η_2 . The components of the convergents are also bounded by $(2d^2)!$, and this is where the extra factor of $(2d^2)!$ come in.

APPENDIX TO CHAPTER 4,

TAKEN FROM [11]

On the reducibility of polynomials

15

Since k_i are integers, $[k_\mu - k_0, k_r - k_0] \in \mathfrak{M}$; thus there are integers s, t such that $k_\mu - k_0 = \xi_1 s, k_r - k_0 = \eta_1 s + \eta_2 t$.

Putting $z_i = (\xi_1 D'_i + \eta_1 D''_i)/D, \lambda_i = \eta_2 D''_i/D$, we get for $i \leq l$

$$k_i - k_0 = z_i s + \lambda_i t$$

and by (39)

$$|z_i| \leq \frac{\xi_1}{|D|} |D'_i| + \frac{|\eta_1|}{D} |D''_i| \leq 2(5c)^{l-1} < (5c)^l,$$

$$|\lambda_i| \leq \frac{\eta_2}{|D|} |D''_i| < (5c)^{l-1}.$$

This completes the proof.

Remark. For a given finite linear set, denote by ϱ the number of rationally independent distances and by ϱ_0 the number of rationally independent distances which appear only once. It follows from the lemma that if $\varrho_0 \leq 2$, then $\varrho \leq 2$. It can easily be found from remark 1 at the end of paper [2] that if $\varrho_0 = 1$ then $\varrho = 1$. The equality $\varrho = \varrho_0$ suggests itself, but I am unable to prove it.

DEFINITION. For a given integral matrix A , $h(A)$ will denote the maximum of absolute values of the elements of A .

LEMMA 5. Let Γ be any given integral matrix 2×2 . For arbitrary positive integers d, n, m there exists an integral matrix

$$M = \begin{bmatrix} v_1 & \mu_1 \\ v_2 & \mu_2 \end{bmatrix}$$

satisfying the conditions:

$$(40) \quad 0 \leq v_i \leq ((2d^2)!)^2, \quad 0 \leq \mu_i \leq ((2d^2)!)^2 \quad (i=1, 2),$$

$$(41) \quad |M| > 0,$$

$$(42) \quad [n, m] = [u, v]M, \quad u, v \text{ integers } \geq 0,$$

and with the following property. If

$$(43) \quad [n, m]\Gamma = [s, t]\Delta,$$

where s, t are integers, Δ is an integral matrix,

$$(44) \quad |\Delta| \neq 0 \quad \text{and} \quad h(\Delta) \leq d,$$

then

$$(45) \quad M\Gamma = T\Delta \quad \text{and} \quad [s, t] = [u, v]T,$$

where T is an integral matrix and

$$(46) \quad h(T) \leq 4d((2d^2)!)^2 h(\Gamma).$$

Proof. Let S be the set of all integral matrices Δ satisfying (43) and (44). Integral vectors $[x, y]$ such that for all $\Delta \in S$ and suitable integers s_Δ, t_Δ , $[x, y]T = [s_\Delta, t_\Delta]\Delta$ form a module, say \mathfrak{M} . By (44) $2d^2 \geq |\Delta| \neq 0$, whence $|\Delta|$ divides $(2d^2)!$.

It follows that $[(2d^2)!, 0] \in \mathfrak{M}$ and $[0, (2d^2)!] \in \mathfrak{M}$. Let ξ_1 be the least positive integer such that, for some η_1 , $[\xi_1, \eta_1] \in \mathfrak{M}$ and let η_2 be the least positive integer such that $[0, \eta_2] \in \mathfrak{M}$. Clearly $[\xi_1, \eta_1]$ and $[0, \eta_2]$ form a basis for \mathfrak{M} and we may assume without loss of generality that $0 \leq \eta_1 < \eta_2$. Hence

$$(47) \quad 0 < \xi_1 \leq (2d^2)!, \quad 0 \leq \eta_1 < \eta_2 \leq (2d^2)!.$$

Let

$$\frac{\eta_1}{\eta_2} = \frac{1}{b_1} - \frac{1}{b_2} - \dots - \frac{1}{b_r}$$

be the expansion of $\frac{\eta_1}{\eta_2}$ into a continued fraction, where b_p are integers > 1 ($1 \leq p \leq r$); if $\eta_1 = 0$ let $r = 0$. Put

$$\begin{aligned} A_{-1} &= -1, \quad B_{-1} = 0; \quad A_0 = 0, \quad B_0 = 1; \\ A_{p+1} &= b_p A_p - A_{p-1}, \quad B_{p+1} = b_p B_p - B_{p-1} \quad (0 \leq p < r). \end{aligned}$$

It follows that the sequences A_p, B_p are increasing and for $p \leq r$

$$(48) \quad A_p B_{p-1} - B_p A_{p-1} = 1,$$

$$(49) \quad 0 \leq A_p \leq \eta_1, \quad 0 < B_p \leq \eta_2,$$

$$(50) \quad A_p/B_p < A_r/B_r = \eta_1/\eta_2.$$

Since m, n are > 0 , we have

$$\frac{\eta_1}{\eta_2} - \frac{m}{n} \cdot \frac{\xi_1}{\eta_2} < \frac{A_r}{B_r}.$$

Let q be the least non-negative integer which can be substituted for r in the last inequality. Assuming $A_{-1}/B_{-1} = -\infty$ we have therefore

$$(51) \quad \frac{A_{q-1}}{B_{q-1}} \leq \frac{\eta_1}{\eta_2} - \frac{m}{n} \cdot \frac{\xi_1}{\eta_2} < \frac{A_q}{B_q}.$$

Let us put

$$M = \begin{bmatrix} \nu_1 & \mu_1 \\ \nu_2 & \mu_2 \end{bmatrix} = \begin{bmatrix} B_q & -A_q \\ B_{q-1} & -A_{q-1} \end{bmatrix} \begin{bmatrix} \xi_1 & \eta_1 \\ 0 & \eta_2 \end{bmatrix} = \begin{bmatrix} B_q \xi_1 & B_q \eta_1 - A_q \eta_2 \\ B_{q-1} \xi_1 & B_{q-1} \eta_1 - A_{q-1} \eta_2 \end{bmatrix}.$$

Inequalities (47), (49) and (50) imply (40). By (48)

$$|M| = \begin{vmatrix} \xi_1 & \eta_1 \\ 0 & \eta_2 \end{vmatrix} = \xi_1 \eta_2 > 0.$$

Moreover, the vectors $[v_1, \mu_1], [v_2, \mu_2]$ form a basis for \mathfrak{M} . Since $[n, m] \in \mathfrak{M}$, there are integers u, v satisfying (42). We have

$$\begin{aligned} [u, v] &= [n, m]M^{-1} = \frac{1}{\xi_1 \eta_2} [n, m] \begin{bmatrix} B_{q-1}\eta_1 - A_{q-1}\eta_2 & -B\eta_1 + A_q\eta_2 \\ -B_{q-1}\xi_1 & B_q\xi_1 \end{bmatrix} \\ &= \frac{1}{\xi_1 \eta_2} [B_{q-1}(n\eta_1 - m\xi_1) - A_{q-1}\eta_2, A_q\eta_2 - B_q(n\eta_1 - m\xi_1)]. \end{aligned}$$

It follows from (51) that $u \geq 0, v \geq 0$. In order to prove the last statement of the lemma suppose that for some integral matrix Δ (43) and (44) hold. Thus $\Delta \in \mathcal{S}$ and since $[v_i, \mu_i] \in \mathfrak{M}$ ($i = 1, 2$) there are integers σ_i, τ_i such that $[v_i, \mu_i]\Gamma = [\sigma_i, \tau_i]\Delta$ ($i = 1, 2$). Putting

$$(52) \quad T = \begin{bmatrix} \sigma_1 & \tau_1 \\ \sigma_2 & \tau_2 \end{bmatrix}$$

we get

$$M\Gamma = T\Delta.$$

On the other hand, (42) and (43) imply

$$(53) \quad [u, v]M\Gamma = [s, t]\Delta.$$

Since $|\Delta| \neq 0$ by (44), we get (45) from (52) and (53). Finally, by (52), (40) and (44)

$$h(T) = h(M\Gamma\Delta^{-1}) \leq 4h(M)h(\Gamma)h(\Delta) \leq 4d((2d^2)!)^2 h(\Gamma).$$

This completes the proof.

LEMMA 6. Let $f(x)$ be an irreducible polynomial not dividing $x^3 - x$ ($\delta > 1$), α, β integers, $\alpha > 0$ or $\beta > 0$. For arbitrary positive integers n, m such that $\alpha n + \beta m > 0$ there exists an integral matrix

$$M = \begin{bmatrix} v_1 & \mu_1 \\ v_2 & \mu_2 \end{bmatrix}$$

satisfying the conditions

$$(54) \quad 0 \leq v_i \leq C(f, \alpha, \beta), \quad 0 \leq \mu_i \leq C(f, \alpha, \beta) \quad (i = 1, 2),$$

$$(55) \quad |M| > 0,$$

$$(56) \quad [n, m] = [u, v]M, \quad u, v \text{ integers } \geq 0,$$

$$(57) \quad \alpha v_i + \beta \mu_i \geq 0 \quad (i = 1, 2),$$

and having the following property:

CHAPTER 5

SUBDIVISION OF 3-DIMENSIONAL LATTICE CONES§5A. Introduction.

Let L be a lattice in R^d . In this chapter I will consider the subdivision problem for 3-dimensional L -cones. As in Chapter 3, a subdivision into basic L -cones of dimension 3 (see Note 3A.1) will be called a basic subdivision, and the terms special subdivision, simplicial subdivision, special complex, etc., will be used similarly.

The argument will hinge on the basic subdivision of special cones (see Definition 3B.9). I shall characterise 3-dimensional special cones in §5B and apply this characterisation to the basic subdivision of special cones in §5C. Then in §5D I shall obtain the main result of this chapter, namely, that any complex S of 3-dimensional L -cones has a basic subdivision T whose base points are indecomposable in S , so that

$$B(T) = I(T) = I(S).$$

The results of §5C, §5D will lead to improvements for the bounds in Theorem 3C.2, in the 3-dimensional case, and these will be given explicitly for a complex of special cones in Theorem 5C.4. A further application of the characterisation of special lattices will be given in §5E.

We conclude this section by settling the question of equivalence of basic subdivisions and basic dissections (see §3A.3.) for the purposes of this chapter.

Proposition 5A.1. Let D be a 3-dimensional L-cone in R^d , and suppose that D_1, D_2 are 3-dimensional basic L-cones contained in D with base points indecomposable in D , and with disjoint relative interiors. Then the intersection $F = D_1 \cap D_2$ is a face of D_1 , and of D_2 .

Proof. By Lemma 1A.1, it follows that no point of F can belong to $\text{relint } D_1$ or to $\text{relint } D_2$. If $F = \{0\}$, or F is an edge common to D_1, D_2 , then the result holds. The remaining two possibilities are

- (i) F is an edge of one of D_1, D_2 lying in a facet of the other, or
- (ii) $F = C_1 \cap C_2$, where C_1, C_2 are facets of D_1, D_2 respectively, $\text{lin } C_1 = \text{lin } C_2$ (and D_1, D_2 lie on opposite sides of F in $\text{lin } D$).

In cases (i), (ii) the result follows easily from Lemma 4A.2 (concerning cones in two dimensions). For example, in case (ii), take $C = D \cap \text{lin } C_1$. Then $C_1, C_2 \subseteq C$. Since D_1, D_2 are basic, so are C_1, C_2 . The base points of C_1, C_2 , being base points of D_1, D_2 respectively, are indecomposable in D , and therefore in C , so the conditions of the lemma are satisfied.

The following result is immediate.

Corollary 5A.2. Suppose that \mathcal{D} is a dissection of a 3-dimensional L-cone C into basic L-cones whose base points are indecomposable in C . Then \mathcal{D} is a subdivision of C .

Since we will be concerned with basic L-cones D in an L-cone C such that the base points of D are indecomposable in C , there is thus no loss of generality in considering subdivisions of C , rather than dissections.

§5B. Characterisation of 3-dimensional special L-cones.

The main result of this section is the characterisation of 3-dimensional special L-cones (Theorem 5B.2). The proof uses Proposition 2.6 from Chapter 2. We begin with two examples.

Example 1. Let $\underline{a}_1 = (1,0,0)$, $\underline{a}_2 = (0,1,0)$, $\underline{a}_3 = (0,0,1)$ in R^3 , and let L be the lattice in R^3 consisting of all points $\frac{1}{2}(\alpha, \beta, \gamma)$, where α, β, γ in Z are all even, or all odd. Then $C = E^+ = c(\underline{a}_1, \underline{a}_2, \underline{a}_3)$ is a special L-cone of index 2. The lattice L has basis

$$\underline{a}_1, \underline{a}_2, \underline{a} = \frac{1}{2}(\underline{a}_1 + \underline{a}_2 + \underline{a}_3) = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}),$$

and \underline{a} is the only non-zero point of $L \cap P(C)$.

Example 2. Let $\underline{u}, \underline{v}, \underline{w}$ be 3 linearly independent points in R^d . Let m, b, c be positive integers such that

$$(1) \quad m > 1, \quad b + c = m, \quad (b, m) = (c, m) = 1.$$

Let L be the lattice generated by $\underline{u}, \underline{v}, \underline{w}$ and the point

$$(2) \quad \underline{a} = \frac{1}{m} (\underline{u} + b\underline{v} + c\underline{w}).$$

That is,

$$(3) \quad L \cap P(C) = \left\{ \frac{1}{m} (r\underline{u} + \overline{r}b \underline{v} + \overline{r}c \underline{w}) \mid 0 \leq r \leq m-1 \right\},$$

where \overline{x} denotes the least non-negative residue of $x \bmod m$. Then the cone

$$C = c(\underline{u}, \underline{v}, \underline{w})$$

is special of index m .

Proof. The index of C is clearly (by Proposition 3B.2) equal to the number of points in $L \cap P(C)$, that is, m . For any point in $L \cap P(C)$, the sum of the coordinates with respect to basis $\underline{u}, \underline{v}, \underline{w}$ is, by (1),

$$(4) \quad \frac{r}{m} + \frac{\overline{rb}}{m} + \frac{\overline{rc}}{m} = \frac{r}{m} + 1 \quad (1 \leq r \leq m-1)$$

$$> 1.$$

Hence C is special of index m .

It will turn out that all special cones of index greater than 1 are of this type. We therefore introduce the following definition.

Definition 5B.1. An L-cone $C = c(\underline{u}, \underline{v}, \underline{w})$ is \underline{u} -special if it is basic or if there exist positive integers m, b, c such that (1), (3) hold. The point (2) will be called the minimal point of C .

Lemma 5B.1. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be an L-cone, and let S be the set of coordinate vectors with respect to $\underline{u}, \underline{v}, \underline{w}$ of the lattice points in the parallelotope of C , that is,

$$(5) \quad S = \{(\lambda, \mu, \nu) \mid \lambda \underline{u} + \mu \underline{v} + \nu \underline{w} \in L \cap P(C)\}.$$

Let m be an integer greater than 1. Then C is \underline{u} -, \underline{v} - or \underline{w} -special of index m if and only if for some a, b, c in \mathbb{Z} ,

$$(6) \quad \begin{cases} S = \left\{ \frac{1}{m} (\overline{ra}, \overline{rb}, \overline{rc}) \mid 0 \leq r \leq m-1 \right\}, \\ 1 \leq a, b, c \leq m-1, \\ \overline{ra} + \overline{rb} + \overline{rc} = m + r \text{ for } r = 1, \dots, m-1. \end{cases}$$

Proof. Suppose that C is \underline{u} -special of index m . Then by Definition 5B.1, (1), (3) hold for suitable b, c . Thus (4), and hence (6) holds with $a = 1$. Similarly, (6) holds if C is \underline{v} - or \underline{w} -special.

Conversely, suppose that (6) holds. Then the index of C , being the number of elements in S , is equal to m .

Also, by Proposition 2.6 with $k = 3$, we have that a, b, c are all relatively prime to m , one is equal to 1, and the other two add up to m . In the case $a = 1$, (1), (3) follow, and C is \underline{u} -special. In the cases $b = 1, c = 1$, C is \underline{v} -special or \underline{w} -special respectively. This completes the proof of the lemma.

We are now in a position to characterise 3-dimensional special cones.

Theorem 5B.2. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be an L-cone. Then C is special if and only if C is \underline{u} -, \underline{v} - or \underline{w} -special.

Proof. The proof of Example 2 shows that if C is \underline{u} -, \underline{v} - or \underline{w} -special then it is special. For the converse, we may suppose that C is special and of index $m > 1$. Let S be defined in (5) of Lemma 5B.1. Then S has m members, all of the form $\frac{1}{m}(\alpha, \beta, \gamma)$, where α, β, γ are non-negative integers less than m . We must show that (6) holds. First we give two preliminary results about S .

(i) Suppose $\frac{1}{m}(\alpha, \beta, \gamma)$ is a non-zero element of S .

Then

$$(7) \quad \begin{cases} m+1 \leq \alpha + \beta + \gamma \leq 2m-1 \\ 1 \leq \alpha, \beta, \gamma \leq m-1. \end{cases}$$

These inequalities are easily deduced from the fact that C is special by considering points such as $(1, 1, 1) - \frac{1}{m}(\alpha, \beta, \gamma)$, $(0, 1, 1) - \frac{1}{m}(\alpha, \beta, \gamma)$ as well as (α, β, γ) .

(ii) Suppose that $\underline{a}_1, \underline{a}_2 \in S$, where $\underline{a}_1 \neq \underline{a}_2$, and for $i = 1, 2$, $\underline{a}_i = \frac{1}{m}(\alpha_i, \beta_i, \gamma_i)$. Then $\alpha_1 + \beta_1 + \gamma_1 \neq \alpha_2 + \beta_2 + \gamma_2$.

Proof of (ii). Suppose to the contrary that $\alpha_1 + \beta_1 + \gamma_1 = \alpha_2 + \beta_2 + \gamma_2$. Since $\alpha_1 \neq \alpha_2$, we may suppose without loss of generality, that $\alpha_1 > \alpha_2$, $\beta_1 \geq \beta_2$, $\gamma_1 < \gamma_2$. Let $\alpha_3 = \alpha_1 - \alpha_2$, $\beta_3 = \beta_1 - \beta_2$, $\gamma_3 = \gamma_1 - \gamma_2 + m$, and let $\underline{a}_3 = \frac{1}{m} (\alpha_3, \beta_3, \gamma_3)$. Then $\underline{a}_3 = \underline{a}_1 - \underline{a}_2 + (0, 0, 1)$ is a non-zero point of S , and $\alpha_3 + \beta_3 + \gamma_3 = m$, contradicting the fact that C is special (the first inequality in (7)).

From (i), (ii) we see that for non-zero elements $\frac{1}{m} (\alpha, \beta, \gamma)$ of S , $\alpha + \beta + \gamma$ takes all the values $m+1, \dots, 2m-1$. Let $\frac{1}{m} (a, b, c)$ be the member of S for which $a + b + c = m+1$. Then by (7), we obtain (6), as required.

Corollary 5B.3. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be a special L-cone. Then C is \underline{u} -special if and only if

$$(8) \quad \mu + \nu = 1 \quad \text{for all } \underline{a} = \lambda \underline{u} + \mu \underline{v} + \nu \underline{w} \text{ in } \dot{L} \cap P(C).$$

Proof. If C is \underline{u} -special, then it is immediate that (8) holds. Conversely, suppose that (8) holds, and suppose that C is \underline{v} -special. Then for all $\underline{a} = \lambda \underline{u} + \mu \underline{v} + \nu \underline{w}$ in $\dot{L} \cap P(C)$, we have $\lambda + \nu = 1$, and by (8), $\mu + \nu = 1$, so $\lambda = \mu$. Hence C is also \underline{u} -special. Similarly if C is \underline{w} -special.

The following corollary is useful in calculations with 3-dimensional L-cones (for example, in checking whether a given complex is special). The denominator q of a non-zero rational number $\lambda = \frac{p}{q}$, $(p, q) = 1$, will be denoted by $\delta(\lambda)$.

Corollary 5B.4. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be an L-cone and let (λ, μ, ν) denote the coordinate vector of \underline{a} with

respect to basis $\underline{u}, \underline{v}, \underline{w}$. If C is \underline{u} -special, then

- (i) $\mu + v \in \mathbb{Z}$ for all $\underline{a} \in \dot{L}$
- (ii) $\delta(\lambda) = \delta(\mu) = \delta(v)$ for all $\underline{a} \in \dot{L}$
- (iii) The index of C is

$$m(C) = \max_{\underline{a} \in \dot{L}} \delta(\lambda)$$

- (iv) $\dot{L} \cap P(C) \subseteq \text{relint } C$.

(Hence, by Theorem 5B.2, (ii)-(iv) hold if C is special.) Conversely, if (i), (ii) hold, then C is \underline{u} -special.

Proof. Let C be \underline{u} -special. Then (i), (ii), (iii) clearly hold. We now prove that (ii) implies (iv). Let $\underline{a} \in \dot{L} \cap P(C)$. Then

$$(9) \quad \underline{a} = \lambda \underline{u} + \mu \underline{v} + \nu \underline{w}, \quad 0 \leq \lambda, \mu, \nu < 1.$$

If $\lambda = 0$, then $\delta(\lambda) = 1$. Hence by (ii), $\delta(\mu) = \delta(\nu) = 1$, and so $\mu = \nu = 0$, so $\underline{a} = \underline{0}$, a contradiction. Thus $\lambda > 0$. Similarly $\mu, \nu > 0$. Thus

$$(10) \quad \lambda, \mu, \nu > 0,$$

and this is equivalent to (iv).

To prove the converse part, let (i), (ii) hold. Let $\underline{a} \in \dot{L} \cap P(C)$. Then (9), (10) hold. By (i), (9), $\mu + \nu$ must be 0 or 1. By (10),

$$\mu + \nu = 1, \quad \lambda + \mu + \nu > 1.$$

Thus C is special, and the condition of Corollary 5B.3 is satisfied. It follows that C is \underline{u} -special.

Corollary 5B.5. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be a special L -cone of index m . Then

$$(11) \quad I(C) = (\dot{L} \cap P(C)) \cup \{\underline{u}, \underline{v}, \underline{w}\}.$$

Thus C has $m+2$ indecomposable points, of which

$\underline{u}, \underline{v}, \underline{w}$ are on the relative boundary of C , and the other $m-1$ points, comprising $\dot{I} \cap P(C)$, are in the relative interior of C .

Proof. Clearly that left hand side is contained in the right hand side. For the reverse inclusion we need only show that

$$(12) \quad \dot{I} \cap P(C) \subseteq I(C).$$

We will use Lemma 3D.1 (iv) with $F = \text{conv}(\underline{u}, \underline{v}, \underline{w})$. By definition of $\underline{n} = \underline{n}(F)$, $\underline{n} \cdot \underline{x} = 1$ for all \underline{x} in F .

Hence, if

$$\underline{a} = \frac{1}{m} (\alpha \underline{u} + \beta \underline{v} + \gamma \underline{w}) \in \dot{I} \cap P(C),$$

then

$$\begin{aligned} \underline{n} \cdot \underline{a} &= \frac{1}{m} (\alpha + \beta + \gamma) \\ &\leq 2 - \frac{1}{m} < 2 \quad \text{by Theorem 5B.2,} \end{aligned}$$

and it follows, from the lemma, that $\underline{a} \in I(C)$. This completes the proof of (12), and hence of (11). The last part of the corollary follows from Corollary 5B.4 (iv) (and Theorem 5B.2).

[The converse does not quite hold. If (12) holds, then from Corollaries 3E.2, 3 we can show that if $\underline{a} \in \dot{I} \cap P(C)$, then

$$\begin{aligned} \underline{a} &= \frac{1}{m} (\alpha \underline{u} + \beta \underline{v} + \gamma \underline{w}), \\ m &\leq \alpha + \beta + \gamma \leq 2m - 1. \end{aligned}$$

But we cannot show that C is special, for, as is easily shown, the left hand bound cannot be improved to $m+1$.]

The following corollary of Corollary 5B.5 relating the number of indecomposable points of a special complex S of 3-dimensional L -cones to the inner index $m_i(S)$ of

S (defined in Definition 3D.1) will be a useful adjunct to the theorem (Theorem 5C.4) on bounds obtained later. If S is a complex of 3-dimensional L-cones, then as in §3C, we let

$$N(S) = \# \{C \mid C \in S\},$$

$$E(S) = \# \{\underline{b} \mid \underline{b} \text{ is a base point of } S\},$$

$$H(S) = \max\{H(\underline{b}) \mid \underline{b} \text{ is a base point of } S\},$$

and if US is 3-dimensional, we also let

$$E_b(S) = \# \{\underline{b} \mid \underline{b} \text{ is a base point of } S \text{ on the relative boundary of } US\},$$

and

$$E_i(S) = \# \{\underline{b} \mid \underline{b} \text{ is a base point of } S \text{ in the relative interior of } US\},$$

so that

$$E(S) = E_i(S) + E_b(S).$$

Corollary 5B.6. Let S be a special complex of 3-dimensional L-cones. Then, in terms of the above notation,

$$(i) \quad \# I(S) = m_i(S) - N(S) + E(S)$$

(ii) If S is a subdivision of an L-cone, or US lies in R^3 and is simply connected, then

$$\# I(S) = m_i(S) + 2 - E_i(S).$$

Proof. (i) By Corollary 5B.5,

$$\begin{aligned} \# I(S) &= \sum_{C \in S} (m(C) - 1) + E(S) \\ &= m_i(S) - N(S) + E(S) \end{aligned}$$

(ii) By Lemma 1C.1.

$$(13) \quad 2E_i(S) + E_b(S) = N(S) + 2.$$

Hence (ii) follows from (i).

The following two results are useful in dealing with \underline{u} -special cones, and are immediate consequences of the definition.

Proposition 5B.7. If $C = c(\underline{u}, \underline{v}, \underline{w})$ is a \underline{u} -special L-cone, then the $m+1$ indecomposable points of $\dot{L} \cap C$ not equal to \underline{u} all lie on the plane through $\underline{v}, \underline{w}$ and parallel to \underline{ou} .

Proposition 5B.8. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be an L-cone. Then C is \underline{u} -special if and only if there is an integer c_1 such that

$$(c_1, m) = 1,$$

$$(\dot{L} \cap P(C)) \cup \{\underline{v}, \underline{w}\} = \{\underline{b}_r \mid r = 0, 1, \dots, m\},$$

$$\underline{b}_r = \frac{1}{m} (\overline{rc_1} \underline{u} + (m-r) \underline{v} + r \underline{w}).$$

§5C. Subdivision of 3-dimensional special L-cones, and complexes.

The main result in this section (Theorem 5C.3) is that every special 3-dimensional complex S of L-cones has a unique basic subdivision T whose base points are indecomposable points of S , so that

$$B(T) = I(T) = I(S)$$

(See Definition 3E.1 and Note 3E.5. The first equality holds because T is basic.)

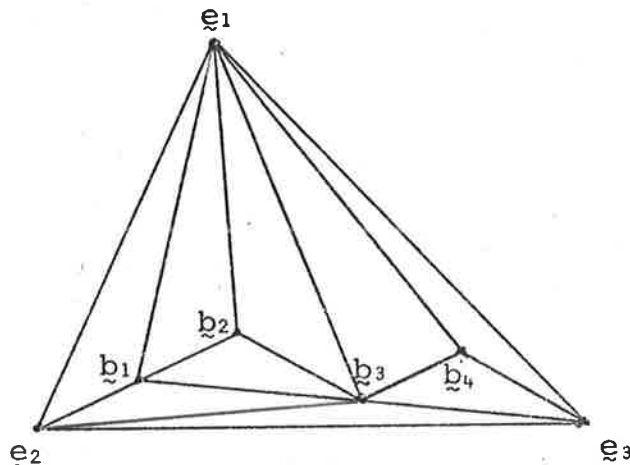
Example 1. Let $C = E^+ = c(\underline{e}_1, \underline{e}_2, \underline{e}_3)$, and let L be given by $(L \cap P(C)) \cup \{\underline{e}_2, \underline{e}_3\} = \{\underline{b}_i \mid i = 0, \dots, 5\}$, where

$$\begin{aligned} \underline{b}_0 &= \frac{1}{5}(0, 5, 0) = \underline{e}_2 \\ \underline{b}_1 &= \frac{1}{5}(2, 4, 1) \\ \underline{b}_2 &= \frac{1}{5}(4, 3, 2) \\ \underline{b}_3 &= \frac{1}{5}(1, 2, 3) \\ \underline{b}_4 &= \frac{1}{5}(3, 1, 4) \\ \underline{b}_5 &= \frac{1}{5}(0, 0, 5) = \underline{e}_3 \end{aligned}$$

(See Proposition 5B.8.)

Then C is \underline{e}_1 -special of index 5, and the diagram below represents a basic subdivision T with indecomposable base points (plane cross-sections being shown).

Subdivision T :



All the 9 ($= 2.5 - 1$) cones have determinant $1/5 = 1/5 |\det(\underline{e}_1, \underline{e}_2, \underline{e}_3)|$, and so are basic. Barycentric subdivision of C with division point \underline{b}_3 produces three \underline{e}_1 -special cones (this could be verified numerically, but will follow from the theory) and further barycentric subdivisions with division points $\underline{b}_2, \underline{b}_1, \underline{b}_4$ results in the subdivision T as shown.

Lemma 5C.1. Let S be a complex of 3-dimensional special cones, and let T be the complex obtained from S by barycentric subdivision with division point \underline{t} in C , where $C \in S$. If $\underline{t} = \underline{t}_0$, the minimal point of C (as a \underline{u} -, \underline{v} -, or \underline{w} -special cone; see Theorem 5B.2 and Definition 5B.1), then T is special, and

$$I(T) = I(S).$$

Proof. Suppose $\underline{t} = \underline{t}_0$, and (using Lemma 5B.1 with $r=1$) let

$$(1) \quad \left\{ \begin{array}{l} C = c(\underline{u}, \underline{v}, \underline{w}), \quad C_{\underline{u}} = c(\underline{t}_0, \underline{v}, \underline{w}), \\ C_{\underline{v}} = c(\underline{u}, \underline{t}_0, \underline{w}), \quad C_{\underline{w}} = c(\underline{u}, \underline{v}, \underline{t}_0), \\ m = m(C) \quad (= \text{index of } C), \\ \underline{t}_0 = \frac{1}{m} (a\underline{u} + b\underline{v} + c\underline{w}), \\ \text{so } a + b + c = m + 1. \end{array} \right.$$

Then T is obtained from S by replacing C by the three cones $C_{\underline{u}}, C_{\underline{v}}, C_{\underline{w}}$, and it is sufficient to show that taking D to be any of these three cones, D is special, and

$$(2) \quad I(D) = D \cap I(C).$$

If $S^*(D)$ is the deleted simplex of D , and

$$\underline{d} = \lambda \underline{u} + \mu \underline{v} + \nu \underline{w} \in L,$$

then

$$\begin{aligned} \lambda + \mu + \nu &< 1 + \frac{1}{m} \quad \text{for } \underline{d} \in L \cap S^*(D), \quad \text{by (1),} \\ &\geq 1 + \frac{2}{m} \quad \text{for } \underline{d} \in \dot{L} \cap C, \\ &\quad \underline{d} \neq \underline{t}_0, \underline{u}, \underline{v}, \underline{w}. \end{aligned}$$

Hence $L \cap S^*(D) = \emptyset$, and D is special.

To prove (2), we note that the right hand side is trivially contained in the left hand side. We will complete the proof of (2) by a counting argument. By Definition 3B.4 and Proposition 3B.1, the indices of the cones $C_{\underline{u}}$, $C_{\underline{v}}$, $C_{\underline{w}}$ are

$$m(C_{\underline{u}}) = a, \quad m(C_{\underline{v}}) = b, \quad m(C_{\underline{w}}) = c.$$

Let

$$a' = \#((C_{\underline{u}} \cap I(C)) \sim \{\underline{t}_0, \underline{v}, \underline{w}\}),$$

and similarly define b' , c' . Now

$$a - 1 = \#(I(C_{\underline{u}}) \sim \{\underline{t}_0, \underline{v}, \underline{w}\})$$

by Corollary 5B.5, and similar results hold for $b-1$, $c-1$.

Thus, by the trivial part of (2), we obtain

$$(3) \quad a' \leq a-1, \quad b' \leq b-1, \quad c' \leq c-1.$$

Thus

$$\begin{aligned} (4) \quad m+2 &= a+b+c+1 \geq a' + b' + c' + 4 \geq \#(I(C)) \\ &= m+2 \quad \text{by Corollary 5B.5.} \end{aligned}$$

Hence equality holds throughout (4) and (3), and therefore (2) holds for $D = C_{\underline{u}}, C_{\underline{v}}, C_{\underline{w}}$. This completes the proof of Lemma 5C.1.

Note. It is easily shown, using Corollary 5B.3 that if C is \underline{u} -special, and $\underline{t} = \underline{t}_0$, then $C_{\underline{u}}$ is basic and $C_{\underline{v}}, C_{\underline{w}}$ are \underline{u} -special.

Note. If $\underline{t} \neq \underline{t}_0$, \underline{u} , \underline{v} , \underline{w} , then it is easily shown that T is not special, and it can be shown from the uniqueness part of Theorem 5C.3 below, that $I(T) \neq I(S)$.

The following lemma partially states the uniqueness part of Theorem 5C.3 and will be used to prove it. It is well illustrated by Example 1 above.

Lemma 5C.2. Let $C = c(\underline{u}, \underline{v}, \underline{w})$ be a \underline{u} -special L-cone of index m , and let

$$\underline{b}_r = \frac{1}{m} (\overline{rc}_1 \underline{u} + (m-r) \underline{v} + r \underline{w})$$

for $r = 0, 1, \dots, m$, and $(c_1, m) = 1$, as in Proposition 5B.8. Let T be a basic subdivision of C for which $I(T) = I(C)$. Then the cones in T containing the point \underline{u} are precisely the $D_r = c(\underline{u}, \underline{b}_{r-1}, \underline{b}_r)$, $r = 1, \dots, m$.

Proof. If $0 \leq s < r \leq m$, then $\det(\underline{u}, \underline{b}_s, \underline{b}_r) = (r-s)/m$, as is easily checked, so that $c(\underline{u}, \underline{b}_s, \underline{b}_r)$ has index $r-s$. The subdivision T must contain some $c(\underline{u}, \underline{b}_0, \underline{b}_r)$, $1 \leq r \leq m$, but this r must be 1 for the cone to be basic. Hence T must contain some $c(\underline{u}, \underline{b}_1, \underline{b}_r)$ such that $\det(\underline{u}, \underline{b}_1, \underline{b}_r) = +1$ and this r must be 2 for the cone to be basic. Proceeding this way, the result follows.

Theorem 5C.3. Let S be a special complex of 3-dimensional L-cones in R^d . Then S has a unique basic subdivision T whose base points are the indecomposable points of S . That is,

$$B(T) = I(T) = I(S)$$

(where the first equality holds because T is basic).

The subdivision T may be obtained from S by repeated barycentric subdivision, where each division point is the minimal point (see Definition 5B.1) of the

relevant cone.

Proof. Existence. This follows from Lemma 5C.1, as the process must terminate (with basic T), because for any subdivision V obtained in the process, $B(V) \subseteq I(S)$, and each barycentric subdivision enlarges $B(V)$ by one point.

Uniqueness. It is sufficient to prove the result in the case when S consists of a single \underline{u} -special cone $C = c(\underline{u}, \underline{v}, \underline{w})$ of index m greater than 1. Suppose that T is a basic subdivision of C such that $I(T) = I(C)$. Since C is special, it follows from Corollary 5B.5 that the facet $c(\underline{v}, \underline{w})$ contains no point of $I(C)$ apart from $\underline{v}, \underline{w}$. Hence for some point \underline{a} in $I(C)$, we have $c(\underline{v}, \underline{w}, \underline{a}) \in T$. By Definition 5B.1, since $c(\underline{v}, \underline{w}, \underline{a})$ is basic, we must have $\underline{a} = \frac{1}{m} (\underline{u} + b\underline{v} + c\underline{w}) = \underline{b}_c$ (in the notation of Lemma 5C.2). Thus the cones $c(\underline{v}, \underline{b}_c)$, $c(\underline{w}, \underline{b}_c)$ lie on the relative boundary of a member of T . By Lemma 5C.2, and in the notation of that lemma, $D_c \in T$, so $c(\underline{u}, \underline{b}_c)$ also lies on the relative boundary of a member of T . Thus, none of these three 2-dimensional cones can contain a point of the relative interior of a member of T (by Lemma 1A.1, as T is a dissection of C). It follows that if $D \in T$, then $\text{relint } D$ lies in one of $C_{\underline{u}}, C_{\underline{v}}, C_{\underline{w}}$ (in the notation (1) of Lemma 5C.1), and so D lies in that cone. Thus $T = T_{\underline{u}} \cup T_{\underline{v}} \cup T_{\underline{w}}$, where $T_{\underline{u}}, T_{\underline{v}}, T_{\underline{w}}$ are basic subdivisions of $C_{\underline{u}}, C_{\underline{v}}, C_{\underline{w}}$ respectively, and $I(T_{\underline{u}}) = I(C_{\underline{u}})$, and similarly for $C_{\underline{v}}, C_{\underline{w}}$. Since $C_{\underline{u}}, C_{\underline{v}}, C_{\underline{w}}$ have smaller indices than C , the proof of the uniqueness result can be carried out by induction. Thus the theorem is proved.

We now give the consequent improvements to Theorem 3C.2 for the particular 3-dimensional case, as in Theorem 5C.3. The notation $N(S)$, $E(S)$, $H(S)$, and so on, will be as in §3C, or as in §5B. The proof of the following theorem will use a combinatorial result (Lemma 1C.1) of Chapter 1.

Theorem 5C.4. Let S be a special complex of 3-dimensional L-cones in R^d , and let T be the basic subdivision of S such that $B(T) = I(T) = I(S)$. Then in the case when $S = \{C\}$, and the index of C is $m(C) = m$, we have

- (i) $N(T) = 2m - 1$
- (ii) $E(T) = m + 2$
- (iii) $H(T) < 2H(C)$, and
 $H(T) = H(C)$ if $C = E^+$ (the non-negative orthant).

In the general case, we have

- (i)' $N(T) = 2m_i(S) - N(S)$
- (ii)' $E(T) = \#I(S) = m_i(S) - N(S) + E(S)$
- (iii)' $H(T) < 2H(S)$

Proof. In (ii)' the first equality is by data, and the second is Corollary 5B.6(i) proved earlier from Corollary 5B.5. Part (ii) is the special case, and is proved similarly, or deduced from (ii)'. Part (i) follows from Corollary 5B.5 and Lemma 1C.1. (The lemma gives §5B (13).) Part (i)' is obtained from result (i) by summation. Part (iii) follows from Definition 5B.1 or from §5B (4). Part (iii)' follows from part (iii).

Notice that from Lemma 1C.1, any enlargement of the set $I(T)$, where T is an arbitrary simplicial subdivision of the special complex S , increases $N(T)$. Hence the unique basic subdivision is most economical, not only from the point of view of number of base points, but also from the point of view of number of cones.

§5D Subdivision of 3-dimensional L-cones and complexes.

In this section we obtain results on general 3-dimensional complexes of L-cones from the results in §5C on special 3-dimensional complexes, and from the results in §3D on minimal special subdivisions which were defined in Definition 3D.1, and exist by Theorem 3D.2.

Theorem 5D.1. Let S be a complex of 3-dimensional L-cones, and let S_1 be any minimal special subdivision of S . Then $I(S_1) = I(S)$.

Proof. The proof is similar to that for Corollary 5B.5. If $\underline{a} = \lambda \underline{u} + \mu \underline{v} + \nu \underline{w} \in I(C)$ for $C = c(\underline{u}, \underline{v}, \underline{w})$ in S_1 , then $\lambda + \mu + \nu < 2$, so $\underline{a} \in I(S)$ by Lemma 3D.1 (iv). Hence $I(S_1) \subseteq I(S)$, and the reverse inclusion is of course trivial.

Theorem 5D.2. Any complex S of 3-dimensional L-cones has a basic subdivision T such that $I(T) = I(S)$.

Proof. From Theorem 5D.1 and Theorem 3D.2 there is a (minimal) special subdivision S_1 of S such that $I(S_1) = I(S)$. Applying the existence part of Theorem 5C.3 to the special complex S_1 we obtain a basic subdivision T of S_1 , and hence of S , such that $I(T) = I(S_1) = I(S)$. This completes the proof of Theorem 5D.2.

(Results on $N(T)$, $E(T)$, $H(T)$ would also follow as in Theorem 5C.4.)

§5E. Applications of §5B.

In this section we will discuss an application of the characterisation of special cones to the characterisation of lattice octahedra in R^3 . Another application is given later in the proof of Theorem 6C.3.

Let a_1, \dots, a_n be linearly independent points in R^n , and consider the convex hull

$$(1) \quad K = \text{conv}(\pm a_1, \dots, \pm a_n).$$

If L is a lattice in R^n which contains a_1, \dots, a_n but contains no other points of K , then K is called a lattice octahedron (with respect to lattice L). The problem of characterising lattice octahedra is the problem of determining for given K , those lattices L for which K is a lattice octahedron with respect to L . Mordell [10] discusses this problem, and gives simple proofs of known results for $n = 3, 4$.

We will now obtain a simple proof of the result for $n = 3$ from our theory of 3-dimensional special cones of §5B.

Theorem 5E.1. Let K , given by (1), be a lattice octahedron with respect to L . Then L has basis a_1, a_2, a_3 , or has basis $a_1, a_2, \frac{1}{2}(a_1 + a_2 + a_3)$.

Proof. The eight cones $c(\pm a_1, \pm a_2, \pm a_3)$ are all special L -cones, and each is thus special with respect to a base point. Let m be their common index. If $m = 1$, then a_1, a_2, a_3 are a basis for L , so suppose $m > 1$. Without loss of generality, we may suppose that $C_1 = c(a_1, a_2, a_3)$ is a_1 -special. Let the minimal point of C_1 be

$$\frac{1}{m} (\underline{a}_1 + h\underline{a}_2 + k\underline{a}_3),$$

$$(h, k) = 1, \quad 1 \leq h, k \leq m-1, \quad h+k = m.$$

Without loss of generality, suppose that $h \leq k$, so that $k \geq \frac{m}{2}$. The cone $C_2 = c(\underline{a}_1, -\underline{a}_2, \underline{a}_3)$ contains the lattice point $\frac{1}{m} (\underline{a}_1 + k(-\underline{a}_2) + k\underline{a}_3)$, and is special. If C_2 is \underline{a}_1 -special, then $k+k = m$ (by the easy part of Corollary 5B.3), so $h = k = 1$, $m = 2$, the minimal point of C_1 is $\frac{1}{2}(\underline{a}_1 + \underline{a}_2 + \underline{a}_3)$, and $\underline{a}_1, \underline{a}_2, \frac{1}{2}(\underline{a}_1 + \underline{a}_2 + \underline{a}_3)$ are a basis for L . If, on the other hand, C_2 is $(-\underline{a}_2)$ -special, or \underline{a}_3 -special, then $1+k = m$. But then the cone $C_3 = c(\underline{a}_1, \underline{a}_2, -\underline{a}_3)$ contains the lattice point $\frac{1}{m} (\underline{a}_1 + \underline{a}_2 + (-\underline{a}_3))$, so since C_3 is special, $m = 2$, and it follows again that $\underline{a}_1, \underline{a}_2, \frac{1}{2}(\underline{a}_1 + \underline{a}_2 + \underline{a}_3)$ are a basis for L . This completes the proof of the theorem.

CHAPTER 6

QUADRATIC FAREY SEQUENCESIntroduction.

This chapter concerns the quadratic analogue of Farey sequences and the analogous determinant properties. The work in this chapter provides some applications of some of the ideas and results of the earlier chapters, and this work was in fact the source of my initial interest in the investigation of the cone subdivision problem.

In §6A I shall introduce quadratic Farey sequences, and then in §6B I shall develop some basic results and notation and outline the plan of the remainder of the chapter.

§6A. The quadratic Farey sequence K_n .

We recall that in §4D we considered the classical Farey sequence F_n of order n and also the sequence G'_n which consists of all reduced fractions of height not exceeding n . We recall that G'_n may be generated from $-1/0$, $0/1$ and $1/0$ by the construction of mediants. From this construction it is immediate that G'_n has the following property:

Determinant property of G'_n . If r/s , r'/s' are any two consecutive members of G'_n ($r/s < r'/s'$), then

$$rs' - r's = -1.$$

We shall often identify a polynomial

$$f(x) = ax^2 + bx + c$$

with the triple

$$f = (a, b, c),$$

and we shall denote its height by

$$|f| = \max(|a|, |b|, |c|).$$

Recall that the height of an algebraic number is defined as the height of its primitive minimal polynomial.

We now define a quadratic analogue of G'_n .

Definition. The quadratic Farey sequence K_n of order n is the sequence of all ^{real} quadratic or rational numbers of height less than or equal to n ordered according to ^{the} natural order of the real numbers.

Let

$$(1) \quad \theta_i \in K_n \quad (i = 1, 2, 3)$$

be consecutive in K_n , and let

$$(2) \quad f_i = (a_i, b_i, c_i) \quad (i = 1, 2, 3)$$

be the corresponding primitive minimal polynomials, so that by definition of K_n ,

$$|f_i| \leq n \quad (i = 1, 2, 3).$$

The f_i are of course unique except for sign. Let

$$(3) \quad D = \begin{vmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{vmatrix}.$$

We shall refer to D as the determinant of the f_i or as the determinant of the θ_i . It is natural to ask whether D has some property analogous to the determinant property of G'_n .

Brown and Mahler [2] introduced the above definition of quadratic Farey sequences, obtained some numerical information (a table of information for $n=5$ is given in [2]), and made some observations and conjectures (without any proofs). In particular, they conjectured that if $\theta_1, \theta_2, \theta_3$ are as in (1), and D is as defined by (2), (3),

then

$$(4) \quad D = 0, 1 \text{ or } -1$$

provided θ_2 is irrational. We will refer to (4) as the determinant property, so that the conjecture was that the determinant property always holds when θ_2 is irrational.

However, it was pointed out by A. Hesterman (private communication) that the result is false for $n=7$, and I have verified Hesterman's calculations. The sequence K_7 has 1311 members (of which 655 are positive) and the determinant property for irrational θ_2 holds everywhere except for one set of irrationals $\theta_1, \theta_2, \theta_3$ in the interval $(0,1)$ (between $1/5$ and $1/4$) and the related sets $-\theta_3, -\theta_2, -\theta_1$; $\theta_3^{-1}, \theta_2^{-1}, \theta_1^{-1}$; $-\theta_1^{-1}, -\theta_2^{-1}, -\theta_3^{-1}$ which lie respectively in the intervals $(-1,0)$, $(1,\infty)$, $(-\infty,1)$. The exception involves the 42nd, 43rd and 44th positive members of K_7 , for which $D = -2$ (when the a_i in (2) are chosen to be positive) :

θ_i	Approximate value	f_i
θ_1	.2265	(7, -6, 1)
θ_2	.2287	(7, 3, -1)
θ_3	.2319	(7, 7, -2)

Despite this, some results about the determinant D of three consecutive elements of K_n will be obtained in this chapter. These will include, in particular, a way of eliminating the most obvious discrepancy when θ_2 is rational (see §6F). Before going further, we need some preliminary results and notation.

§6B. Preliminary results for K_n : elements in $(r/s, r'/s')$.

In this chapter we will be considering three consecutive members $\theta_1, \theta_2, \theta_3$ of K_n such that θ_2 is irrational, and as far as the determinant property in K_n is concerned, we will be able to suppose that the θ_i all lie in some closed Farey interval $[r/s, r'/s']$, that is an interval $[r/s, r'/s']$, where $r/s, r'/s'$ are consecutive elements of G'_n . For if $\theta_i \geq n$ for $i = 1, 2, 3$, we may consider the θ_i^{-1} , which lie in $[0, 1/n]$, and for $\theta_i \leq -n$ we may consider $[-1/n, 0]$. (We may also suppose without loss of generality, when convenient, that the θ_i are non-negative.)

Lemma 6B.1. Let $r/s, r'/s'$ be consecutive elements of G'_n .

- (i) The elements of K_n in the open interval $(r/s, r'/s')$ are all irrational.
- (ii) If f is a quadratic polynomial such that $|f| \leq n$, and $f(\theta) = 0$, where $\theta \in (r/s, r'/s')$, then $\theta \in K_n$, and f is an integral multiple of the primitive minimal polynomial of θ .

Proof. (i) This is immediate from the definitions.
 (ii) It is easily checked that any integral linear factor of a reducible integral polynomial of height not exceeding n itself has height not exceeding n . Result (ii) easily follows from this and from (i).

Lemma 6B.2. Let $r/s, r'/s'$ be consecutive elements of G'_n , let $\theta \in K_n \cap (r/s, r'/s')$, let f be the primitive minimal polynomial of θ , and let θ' be the conjugate of θ . Then

- (i) $\theta' \notin (r/s, r'/s')$, and
(ii) we may "normalise" f by multiplying by ± 1 so that
either

$$(1) \quad \begin{cases} f(x) < 0 & \text{for } x \text{ in } (r/s, \theta) \\ f(x) > 0 & \text{for } x \text{ in } (\theta, r'/s'), \end{cases}$$

or (1) holds with the inequality symbols interchanged.

Proof. (i) (Notice that Lemma 6B.1(i) implies that f is quadratic.) Suppose to the contrary that θ, θ' lie in the same Farey interval $(r/s, r'/s')$. Then the conjugate elements $-\theta', -\theta$ lie in the same Farey interval $(-r'/s', -r/s)$, and similarly for $\theta'^{-1}, \theta^{-1}$ and for $-\theta^{-1}, -\theta'^{-1}$. Thus we may suppose that

$$\begin{aligned} 0 &\leq r/s < \theta, \theta' < r'/s' \leq 1, \\ s &\geq 1, s' \geq 1, r \geq 0, r' \geq 1, \\ s + s' &\geq n + 1. \end{aligned}$$

Then

$$\frac{r'}{s'} - \frac{r}{s} = \frac{1}{ss'} \leq \frac{1}{s+s'-1} \leq \frac{1}{n},$$

contradicting the fact that

$$1\theta - \theta'1 = \frac{\sqrt{b^2 - 4ac}}{|a|} > \frac{1}{n}.$$

- (ii) This is immediate from (i).

Remark. The lemma gives a one-one correspondence between the θ in K_n in a given Farey interval, and their quadratics f . Hence we will sometimes refer to the order of quadratics f , when we mean the order of their zeros in $(r/s, r'/s')$.

If $f(x) = ax^2 + bx + c$ is a quadratic, it will sometimes be convenient to consider the quadratic form

$$(2) \quad f(x, y) = ax^2 + bxy + cy^2 = y^2 f(x/y).$$

Definition 6B.1. A quadratic f given by (2) has form Δ at r/s (where r, s are relatively prime) if

$$f(r, s) = \Delta.$$

An element θ in K_n has form Δ if its primitive minimal polynomial has form $\pm\Delta$.

The plan of the remaining sections of this chapter is as follows. In §6C I shall derive some general results concerning pairs and triples of consecutive elements of K_n . In §6D I shall specify the elements of K_n in $[0, 1/n]$ completely and show that the determinant property holds in this interval. In §6E I shall obtain corresponding results for $\left[\frac{n-1}{n}, 1\right]$. In §6F I shall show how to obtain the smallest element of K_n in $(r/s, r'/s')$ for consecutive elements $r/s, r'/s'$ of G'_n , I will show that it has form 1 at r/s , and using this, I will show how to eliminate the discrepancy in the determinant when the middle element is rational. In §6G I will give results on the second smallest element of K_n in $(r/s, r'/s')$, without going into all the details of the proof. In §6H I shall investigate the set S_1 of elements of K_n in $(r/s, r'/s')$ with form 1 at r/s , and show that the determinant property holds on S_1 (that is, for triples of consecutive elements in S_1).

We remark that we may view many of the problems concerning the elements of K_n in terms of lattice points (a point in \mathbb{Z}^3 corresponds to a polynomial $f = (a, b, c)$) and lattice cones. Thus we will use a result on 3-dimensional special cones from Chapter 5 in §6C, and we will use element-

ary results on 2-dimensional lattice cones in §6H.

§6C: General results for 2 or 3 consecutive members of K_n .

Two consecutive quadratics.

Theorem 6C.1. Let $r/s, r'/s'$ be consecutive elements of G'_n , and let

$$f_i = (a_i, b_i, c_i) \quad (i = 1, 2)$$

be the primitive minimal polynomials of two consecutive elements θ_1, θ_2 of K_n in $[r/s, r'/s']$. Then

$$(1) \quad \ell(f_1, f_2) = Z^3 \cap \text{lin}(f_1, f_2),$$

that is, f_1, f_2 are a basis for the 2-dimensional sublattice of Z^3 in the plane of f_1, f_2 .

Proof. In the terminology of §3A.1, we must prove that the Z^3 -cone

$$C = c(f_1, f_2)$$

is basic. We may suppose that $\theta_1 < \theta_2$ and that f_1, f_2 are normalised according to §6B(1), so that

$$f_1(\theta_1) = 0, \quad f_1(\theta_2) > 0, \quad |f_1| \leq n,$$

$$f_2(\theta_1) < 0, \quad f_2(\theta_2) = 0, \quad |f_2| \leq n.$$

Suppose that C is not basic. Then by Lemma 3B.5, C is not special, so there exist λ_1, λ_2 such that

$$f_3 = \lambda_1 f_1 + \lambda_2 f_2 \in Z^3 \cap C$$

$$\lambda_1 + \lambda_2 \leq 1, \quad 0 < \lambda_i < 1 \quad \text{for } i = 1, 2.$$

It follows that

$$|f_3| \leq n, \quad f_3(\theta_1) < 0, \quad f_3(\theta_2) > 0,$$

so

$$f_3(\theta) = 0 \quad \text{for some } \theta \in (\theta_1, \theta_2).$$

By Lemma 6B.1(ii), $\theta \in K_n$, contradicting the consecutiveness of θ_1, θ_2 . Thus C is basic.

Corollary 6C.2. Under the conditions of Theorem 6C.1, the 2×2 minors of the matrix

$$\begin{pmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \end{pmatrix}$$

have no common factor.

Proof. By Lemma 2 of Chapter I of Cassels [3], (1) is equivalent to the condition on the minors.

Three consecutive quadratics.

Theorem 6C.3. Let $r/s, r'/s'$ be consecutive elements of G'_n , let $\theta_1, \theta_2, \theta_3$ ($\theta_1 < \theta_2 < \theta_3$) be three consecutive members of K_n in $[r/s, r'/s']$, and let

$$f_i = (a_i, b_i, c_i) \quad (i = 1, 2, 3)$$

be their primitive minimal polynomials. Let

$$D = \det(f_1, f_2, f_3) = \begin{vmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{vmatrix}.$$

Then

$$|D| \leq n.$$

Proof. We may suppose that the f_i are normalised according to §6B(1) so that

$$\begin{aligned} f_1(\theta_1) &= 0, & f_1(\theta_2) &> 0, & f_1(\theta_3) &> 0 \\ f_2(\theta_1) &< 0, & f_2(\theta_2) &= 0, & f_2(\theta_3) &> 0 \\ f_3(\theta_1) &< 0, & f_3(\theta_2) &< 0, & f_3(\theta_3) &= 0. \end{aligned}$$

Suppose that $|D| > n$. Then the f_i are linearly independent, and the Z^3 -cone

$$C = c(f_1, f_2, f_3)$$

has index $|D| > 1$.

We will show that C is special. Suppose not. Then there are numbers $\lambda_1, \lambda_2, \lambda_3$ in $[0, 1)$ no two of which are zero, such that

$$\lambda_1 + \lambda_2 + \lambda_3 \leq 1$$

and such that

$$f_4 = \lambda_1 f_1 + \lambda_2 f_2 + \lambda_3 f_3 \in \mathbb{Z}^3.$$

It follows that

$$|f_4| \leq n, \quad f_4(\theta_1) < 0, \quad f_4(\theta_3) > 0,$$

so

$$f_4(\theta) = 0 \quad \text{for some } \theta \in (\theta_1, \theta_3).$$

By Lemma 6B.1(ii), $\theta \in K_n$, and f_4 is an integral multiple of the primitive minimal polynomial of θ . Since f_4 is not a multiple of f_2 , $\theta \neq \theta_2$. This contradicts the consecutiveness of $\theta_1, \theta_2, \theta_3$. Thus we have proved that C is special.

By Theorem 5B.2, C is f_1 -, f_2 - or f_3 -special. Hence, by Lemma 5B.1 (taking $r=1$ in §5B(6)), there exist positive μ_1, μ_2, μ_3 such that

$$f_5 = \mu_1 f_1 + \mu_2 f_2 + \mu_3 f_3 \in \mathbb{Z}^3,$$

$$\mu_1 + \mu_2 + \mu_3 = 1 + \frac{1}{|D|}.$$

Thus

$$|f_5| \leq n \left(1 + \frac{1}{|D|} \right) < n+1,$$

and so $|f_5| \leq n$.

As before, this contradicts the consecutiveness of $\theta_1, \theta_2, \theta_3$. Thus we have proved that $|D| \leq n$.

In the next two sections we will study two rational Farey intervals within which the determinant D of 3 consecutive elements of K_n always satisfies

$$|D| \leq 1.$$

§6D Elements of K_n in $[0, 1/n]$.

The elements of K_n in $(0, 1/n)$ may be obtained by taking reciprocals of elements of K_n greater than n , so suppose that θ is a zero greater than n of the integral quadratic

$$f(x) = ax^2 + bx + c \quad (a > 0, |f| \leq n).$$

Then

$$\theta = \frac{-b}{2a} \pm \sqrt{\left(\frac{b}{2a}\right)^2 - \frac{c}{a}}.$$

We may then deduce the following results in the order given:

$$\left|\frac{b}{2a}\right| \leq \frac{n}{2},$$

the $+$ sign must hold,

$$\left|\sqrt{\left(\frac{b}{2a}\right)^2 - \frac{c}{a}}\right| > \frac{n}{2}.$$

$$c < 0,$$

$$a = 1, \quad (a \geq 2 \Rightarrow \theta \leq \frac{n}{4} + \frac{1}{4} \sqrt{n^2 + 8n} \left\{ \begin{array}{l} < \frac{n}{2} + 1 \leq n \text{ if } n \geq 2 \\ = n \text{ if } n = 1 \end{array} \right.$$

$$b = -n \quad (b \geq -n + 1 \Rightarrow \theta \leq \frac{n-1}{2} + \frac{1}{2} \sqrt{(n-1)^2 + 4n} = n),$$

$$f = (1, -n, c),$$

$$\theta = \frac{n + \sqrt{n^2 + 4|c|}}{2}.$$

Hence the possible quadratics $f(x)$ are, in order of ascending zero greater than n ,

$$(1, -n, -1)$$

$$(1, -n, -2)$$

$$\dots$$

$$(1, -n, -n).$$

For each of these quadratics, the positive zero lies between n and $n+1$ exclusive. We may now deduce (using Lemma 6B.1(ii)) that the elements of K_n between 0 and

$1/n$ have primitive minimal polynomials, which when arranged in ascending order, with the reducible polynomials $(1,0,0)$, $(n,-1,0)$ inserted, give the table:

<u>quadratic</u>			<u>determinant</u>
1	0	0	
0	1	0	
n	n	-1	-1
n-1	n	-1	1
.	.	.	0
2	n	-1	:
1	n	-1	:
n	-1	0	0
n	-1	0	-1
0	n	-1	1

We also have

$$\begin{vmatrix} 2 & n & -1 \\ 1 & n & -1 \\ 0 & n & -1 \end{vmatrix} = 0.$$

Thus we have verified that $|D| \leq 1$ for any three consecutive members of K_n in $[0/1, 1/n]$, as well as verifying this determinant property in the slightly modified form of the above table. The first and last figures in the last column of the table illustrate a general result proved later, in §6F, that the least and greatest elements of K_n in $(r/s, r'/s')$, f_1, f_2 , say, when normalised according to 6B(1), satisfy

$$f_1(r, s) = -1, \quad f_2(r', s') = 1.$$

§6E Elements of K_n in $\left[\frac{n-1}{n}, \frac{1}{1}\right]$.

By Lemma 6B.1(ii) the elements of K_n in $\left(\frac{n-1}{n}, 1\right)$ correspond to the primitive quadratics

(1) $f(x) = ax^2 + bx + c,$

for which

$$(2) \quad |f| \leq n,$$

$$(3) \quad \begin{cases} f(n-1, n) = n^2 f\left(\frac{n-1}{n}\right) = a(n-1)^2 + b(n-1)n + cn^2 < 0, \\ \Delta = f(1, 1) = f(1) = a + b + c > 0. \end{cases}$$

Taylor's Theorem gives

$$(4) \quad \begin{cases} f(x) = f(1) + (x-1)\left(f'(1) + \frac{x-1}{2} f''(1)\right) \\ = \Delta + (x-1)\left(2a + b + \frac{x-1}{2} 2a\right). \end{cases}$$

Since

$$\frac{-1}{2n} < \frac{x-1}{2} < 0 \quad \text{for } x \in \left(\frac{n-1}{n}, 1\right),$$

it follows from (1), (2) that the graphs of the $f(x)$, for fixed Δ , do not cross in $\left(\frac{n-1}{n}, 1\right)$.

We now show that when (3) is satisfied too, then $\Delta = 1$, with at most 2 exceptions where $\Delta = 2$. The determinant

$$\begin{vmatrix} n & -(n-1) & 0 \\ 1 & -1 & 0 \\ 0 & 1 & -1 \end{vmatrix} = 1.$$

Hence an arbitrary integer triple is an integral linear combination of the rows. The quadratics corresponding to these rows have, respectively, $\Delta = 1, 0, 0$. Thus an arbitrary integer quadratic may be written as

$$(5) \quad f = (a, b, c) = \Delta(n, -n+1, 0) + \alpha(1, -1, 0) + \beta(0, 1, -1),$$

where α, β are integral, and $\Delta = a + b + c$. It follows that

$$(6) \quad f(n-1, n) = -\alpha(n-1) - \beta n.$$

From (5), (6) it follows that (2), (3) can only hold if $\Delta \leq 2$.

The case $\Delta = 2$ yields at most two primitive quadratics:

$$f_0 = (n-1, 3, -n) \quad (\text{when } n \geq 3)$$

for which

$$f_0(n-1, n) = -1, \quad \frac{1}{2}f_0(1) = 1$$

and

$$g = (n, 2, -n) \quad (\text{when } n \geq 3 \text{ and is odd})$$

for which

$$g(n-1, n) = -n, \quad \frac{1}{2}g(1) = 1.$$

(When n is even, g is not primitive.)

The case $\Delta = 1$ yields $\frac{n^2+3n-2}{2}$ primitive quadratics,

with zeros in $\left(\frac{n-1}{n}, 1\right)$. Their order is determined by (4).

For $n=1$, there is only the quadratic $(1, 1, -1)$ and $n=1$ is already covered by §6D. For $n \geq 2$, we have:

	<u>a</u>	<u>b</u>	<u>c</u>	<u>α</u>	<u>β</u>	<u>$f(n-1, n)$</u>
<u>$2a + b = n + 1$</u>	$n-1$	$-n+3$	-1	-1	1	-1
	$n-2$	$-n+5$	-2	-2	2	-2

	1	$n-1$	$-n+1$	$-n+1$	$n-1$	$-n+1$
<u>$2a + b = n + 2$</u>	n	$-n+2$	-1	0	1	$-n$
	$n-1$	$-n+4$	-2	-1	2	$-n-1$

	1	n	$-n$	$-n+1$	n	$-2n+1$
<u>$2a + b = n + 3$</u>	n	$-n+3$	-2	0	2	$-2n$
	$n-1$	$-n+5$	-3	-1	3	$-2n-1$

	2	$n-1$	$-n$	$-n+2$	n	$-3n+2$
<u>$2a + b = n + 4$</u>	n	$-n+4$	-3	0	3	$-3n$
	$n-1$	$-n+6$	-4	-1	4	$-3n-1$

	3	$n-2$	$-n$	$-n+3$	n	$-4n+3$

	<u>a</u>	<u>b</u>	<u>c</u>	<u>α</u>	<u>β</u>	<u>$f(n-1, n)$</u>
<u>$2a + b = n + 5$</u>
.
<u>$2a + b = 2n$</u>	n	0	$-n+1$	0	$n-1$	$-n(n-1)$
	$n-1$	2	$-n$	-1	n	$-n(n-1)-1$
<u>$2a + b = 2n + 1$</u>	n	1	$-n$	0	n	$-n^2$

Notice that

$$a + b + c = 1$$

$$\alpha + \beta = 2a + b - (n+1).$$

The general rule for the ordering is as follows:

Let the quadratics be

$$f_i = (a_i, b_i, c_i) \quad \left(i = 1, \dots, N, \quad N = \frac{n^2+3n-2}{2} \right),$$

and let the corresponding zeros in $\left(\frac{n-1}{n}, 1 \right)$ be

$$\theta_1 < \theta_2 < \dots < \theta_N.$$

Then (from (4))

$$\theta_i < \theta_j \Leftrightarrow 2a_i + b_i < 2a_j + b_j \quad \text{or}$$

$$2a_i + b_i = 2a_j + b_j \quad \text{and} \quad a_i > a_j.$$

Let us insert into this table the quadratics for which

$\Delta = 2$. These were obtained above, and are:

$$f_0 = (a_0, b_0, c_0) = (n-1, 3, -n) \quad (n \geq 3)$$

$$f_0(\theta_0) = 0, \quad \frac{n-1}{n} < \theta_0 < 1, \quad \text{and}$$

$$g = (a', b', c') = (n, 2, -n) \quad (n \geq 3, \quad n \text{ odd})$$

$$g(\phi) = 0, \quad \frac{n-1}{n} < \phi < 1.$$

Consider the graph of $\frac{1}{2}f_0$. Since

$$\frac{1}{2}(a_0 + b_0 + c_0) = 1,$$

$$\frac{1}{2}(2a_0 + b_0) = n + \frac{1}{2}, \quad \text{and for } i \geq 1,$$

$$\frac{x-1}{2} (\frac{1}{2}a_0 - a_i) > -\frac{1}{2} \quad \text{in } \left(\frac{n-1}{n}, 1 \right),$$

it follows from (4) that the graph of $\frac{1}{2}f_0$ lies above that of any f_i ($i \geq 1$), so that

$$\theta_0 \leq \theta_i \quad \text{for } i \geq 1.$$

As for $\frac{1}{2}g$, we have

$$\frac{1}{2}(a' + b' + c') = 1,$$

$$\frac{1}{2}(2a' + b') = n + 1,$$

so the graph of $\frac{1}{2}g$ lies in $\left(\frac{n-1}{n}, 1\right)$ between the graphs of two quadratics in the first "block" of $n-1$ quadratics for which $2a + b = n + 1$, namely between

$$\left(\frac{n}{2} + \frac{1}{2}, 0, \frac{-n}{2} + \frac{1}{2}\right) \quad \text{and} \quad \left(\frac{n}{2} - \frac{1}{2}, 2, \frac{-n}{2} - \frac{1}{2}\right)$$

and ϕ lies between the corresponding zeros. Thus ϕ lies strictly between the first and last quadratics in that block.

We are now able to calculate the determinants of 3 consecutive quadratic irrationals in $\left(\frac{n-1}{n}, 1\right)$. Let us initially insert the "normalised" quadratics $\frac{1}{2}f_0$, $\frac{1}{2}g$ into the appropriate positions. Then the quadratic polynomials fall into "blocks" according to the value of $2a + b$. When $n = 2$, there are 3 blocks, but when $n \geq 3$, there are $n + 2$ blocks, the extra block being formed by $\frac{1}{2}f_0$, for which

$$2a + b = \frac{1}{2}(2a_0 + b_0) = n + \frac{1}{2}.$$

Suppose that

$$h_i = (p_i, q_i, r_i) \quad (i = 1, 2, 3)$$

are consecutive, and let

$$D = \begin{vmatrix} p_1 & q_1 & r_1 \\ p_2 & q_2 & r_2 \\ p_3 & q_3 & r_3 \end{vmatrix}.$$

Since $\Delta = p + q + r = 1$ for normalised quadratics,

$$D = \begin{vmatrix} p_1 - p_2 & q_1 - q_2 \\ p_2 - p_3 & q_2 - q_3 \end{vmatrix}.$$

(i) If the h_i are from the same block, then they differ by multiples of $(-1, 2, -1)$, so $D = 0$.

If the h_i are not all from the same block, then two equations must come from the same block, and the other from an adjoining block, since all blocks have at least two equations, except for the last block ($2a + b = 2n + 1$) and the first block ($2a + b = n + \frac{1}{2}$ if $n \geq 3$, and $2a + b = n + 1$ if $n = 2$) which have one equation each.

Suppose that h_1, h_2 are from one block, and h_3 from the next. Then

(ii) if $h_1 \neq \frac{1}{2}g$, then

$$D = \begin{vmatrix} 1 & -2 \\ k & -2k-1 \end{vmatrix} = -1 \quad (k \in \mathbb{Z}),$$

(iii) if $h_1 = \frac{1}{2}g$, then

$$D = \begin{vmatrix} \frac{1}{2} & -1 \\ k & -2k-1 \end{vmatrix} = -\frac{1}{2} \quad (k \in \mathbb{Z}).$$

Now suppose that h_1 is from one block, and h_2, h_3 are from the next. Then

(iv) if $h_1 \neq \frac{1}{2}f_0$, then

$$D = \begin{vmatrix} k & -2k-1 \\ 1 & -2 \end{vmatrix} = 1 \quad (k \in \mathbb{Z}),$$

(v) if $h_1 = \frac{1}{2}f_0$, and $h_3 \neq \frac{1}{2}g$, then

$$D = \begin{vmatrix} k & -2k-\frac{1}{2} \\ 1 & -2 \end{vmatrix} = \frac{1}{2} \quad (k \in \mathbb{Z}),$$

(vi) if $h_1 = \frac{1}{2}f_0$, and $h_3 = \frac{1}{2}g$ (this can only happen if $n = 3$), then

$$D = \begin{vmatrix} k & -2k-\frac{1}{2} \\ \frac{1}{2} & -1 \end{vmatrix} = \frac{1}{4} \quad (\text{in fact } k = 1).$$

Thus, in all cases, on reverting where necessary from the normalised forms $\frac{1}{2}f_0$, $\frac{1}{2}g$ to the primitive forms f_0 , g , the determinant of 3 consecutive quadratic irrationals in $\left(\frac{n-1}{n}, 1\right)$ is $0, \pm 1$.

It may be checked that if we add one or both of $(n, -n+1, 0)$, $(0, n, -n+1)$ to the beginning of the table, and if we add one or both of $(1, -1, 0)$, $(0, 1, -1)$ to the end of the table, then the additional determinants obtained are all equal to $0, \pm 1$. It follows in particular that the determinant property holds for K_n in the closed interval $\left[\frac{n-1}{n}, 1\right]$.

§6F. The smallest element of K_n in $(r/s, r'/s')$.

Throughout this section we let $r/s, r'/s'$ be consecutive elements of G'_n ($r/s < r'/s'$).

We will show that the least element of K_n in $(r/s, r'/s')$ has form 1 at r/s . (See Definition 6B.1.) Notice that results about the least element of K_n in $(r/s, r'/s')$ may be converted into results about the greatest element of K_n in the Farey interval $(s'/r', s/r)$ by considering the mapping $\theta \rightarrow \theta^{-1}$ ($\theta \in K_n$). For example, the result stated above implies that the greatest element of K_n in a Farey interval $(r/s, r'/s')$ has form 1 at r'/s' .

Also, results for $r/s \leq 0$ may be deduced from results for $r/s \geq 0$ by considering the mapping $\theta \rightarrow -\theta$ ($\theta \in K_n$). Thus we will often be able to assume that $r/s \geq 0$ without loss of generality. Notice that the case $r/s = 0$ has been considered fully in §6D.

Incidentally, it is an easy matter to show that there do exist elements of K_n in $(r/s, r'/s')$. For example, there exists such an element corresponding to one of the quadratics

$$(s, -r, 0) \pm (0, s', -r').$$

This is easily checked using the result

$$(1) \quad r's - s'r = 1.$$

We prove the following stronger assertion.

Lemma 6F.1. Let $r/s, r'/s'$ be consecutive elements of G'_n . Then there is an element of K_n in $(r/s, r'/s')$ having form 1 at r/s .

Proof. When $r/s = 0$, the quadratic $(-n, -n, 1)$ determines an element of form 1 at 0 in $(0, 1/n)$. Thus we may suppose that $r/s \neq 0$, and without loss of generality, $r/s > 0$. Thus $r \geq 1, s \geq 1$.

Since $\text{g.c.d.}(r, s) = 1$, the equation

$$(2) \quad f(r, s) = ar^2 + brs + cs^2 = 1$$

has an integral solution

$$(3) \quad f = f_0 = (a_0, b_0, c_0), \quad f_0(r, s) = 1.$$

(In fact, $f = (s'^2, -2r's', r'^2)$ is a solution by (1).)

As (3) can be written

$$(4) \quad \begin{vmatrix} a_0 & b_0 & c_0 \\ s & -r & 0 \\ 0 & s & -r \end{vmatrix} = 1,$$

it follows that the rows of the matrix are a basis for \mathbb{Z}^3 , and hence the general solution of (2) is

$$(5) \quad f = f_0 + \alpha(s, -r, 0) + \beta(0, s, -r) \quad (\alpha, \beta \in \mathbb{Z}).$$

Let $f = f_1$ be given by (3), (5) and

$$(6) \quad f = f_1 = (a_1, b_1, c_1) \Leftrightarrow \alpha = \left[\frac{-a_0}{s} \right], \quad \beta = \left[\frac{c_0}{r} \right].$$

Then, denoting the fractional part of x by $\{x\}$, it follows from (5), (6) that

$$(7) \quad a_1 = -s \left\{ \frac{-a_0}{s} \right\}, \quad c_1 = r \left\{ \frac{c_0}{r} \right\}.$$

By using (5), (6) and (3) it is also easily checked that

$$(8) \quad b_1 = r \left\{ \frac{-a_0}{s} \right\} - s \left\{ \frac{c_0}{r} \right\} + \frac{1}{rs}.$$

From (7), (8), $|f_1|$ is at most $\max(r, s)$. However, a further adjustment may be necessary to the solution (6) of (2), since $f_1(r', s')$ might not be negative, and by Lemma 6B.1(ii') this is a necessary and sufficient condition for f_1 to give an element of K_n .

From (5), it follows that

$$(9) \quad f(r', s') = f_0(r', s') + \alpha r' + \beta s'.$$

Hence from (6), it may be verified, using (1), (3), that

$$(10) \quad f_1(r', s') = -r' \left\{ \frac{-a_0}{s} \right\} - s' \left\{ \frac{c_0}{r} \right\} + \frac{r's'}{rs}.$$

Let r'', s'' be the integers determined uniquely by

$$(11) \quad r''s - rs'' = 1, \quad 0 < r'' \leq r, \quad 0 \leq s'' < s.$$

(The first and third conditions imply the second, as $r > 0$.)

Then

$$\begin{bmatrix} r' \\ s' \end{bmatrix} = \begin{bmatrix} r'' \\ s'' \end{bmatrix} + k \begin{bmatrix} r \\ s \end{bmatrix}, \quad k \in \mathbb{Z}, \quad k \geq 0$$

(and $k > 0$ if $s'' = 0$). We have

$$(12) \quad \left\lfloor \frac{s'}{s} \right\rfloor = k; \quad \left\lfloor \frac{r'}{r} \right\rfloor = k \quad \text{if } r > 1, \quad \text{and}$$

$$(13) \quad \frac{r'}{r} = k + 1 \quad \text{if } r = 1.$$

Since $r/s, r'/s'$ are consecutive in G'_n , k may also be described as the largest non-negative integer such that (for given r, s, r'', s'', n)

$$(14) \quad r' = kr + r'' \leq n, \quad s' = ks + s'' \leq n$$

both hold.

In the remainder of the proof, there are three cases to consider.

Case $r \geq 2, s \geq 2$. In this case, (10) gives

$$f_1(r', s') \leq -\frac{r'}{s} - \frac{s'}{r} + \frac{r's'}{rs}.$$

Since, by (1),

$$\frac{r'}{r} = \frac{s'}{s} + \frac{1}{rs} > \frac{s'}{s},$$

it follows that

$$\begin{aligned} f_1(r', s') &< \frac{-s'}{s} \left(\frac{r}{s} + \frac{s}{r} - \frac{r'}{r} \right), \quad \text{so} \\ (15) \quad f_1(r', s') &< \frac{s'}{s} \left(\frac{r'}{r} - 2 \right). \end{aligned}$$

Now define $f = f_2$ by (5), substituting f_1, α_1, β_1 for f_0, α, β respectively. Then (9) gives

$$(16) \quad f_2(r', s') = f_1(r', s') + \alpha_1 r' + \beta_1 s'.$$

Choose

$$(\alpha_1, \beta_1) = \begin{cases} (0, 0) & \text{if } k = 0, 1 \\ [-(k-1), -(k-1)] & \text{if } k \geq 2. \end{cases}$$

Then from (12), (15), (16), $f_2(r', s') < 0$, and from (7), (8), (14), $|f_2| \leq n$. Thus f_2 is the primitive (all polynomials of form 1 are primitive) minimal polynomial of an element of K_n in $(r/s, r'/s')$ having form 1 at r/s .

Case $r = 1$. Here, $r'' = 1, s'' = s - 1$. Also,

$a_0 \equiv 1(s)$. Thus

$$(17) \quad f_1 = (-s+1, 1, 0) = (-s'', r'', 0).$$

From (10), or directly from (17),

$$f_1(r', s') = kr'.$$

Define f_2 as before, choosing

$$(\alpha_1, \beta_1) = \begin{cases} (0, -1) & \text{if } k = 0 \\ (-k, -k) & \text{if } k \geq 1. \end{cases}$$

Then f_2 gives an element of K_n as required.

Case $s = 1$. Here, we proceed as in the previous case to get

$$f_1 = (0, 0, 1) = (-s'', s'', r''),$$

$$f_1(r', s') = s'^2 = ks',$$

choosing

$$(\alpha_1, \beta_1) = \begin{cases} (-1, 0) & \text{if } k = 0 \\ (-k, -k) & \text{if } k \geq 1. \end{cases}$$

Then f_2 gives the required element of K_n . This completes the proof of Lemma 6F.1.

The following corollary is immediate.

Corollary 6F.2. Let $r/s, r'/s'$ be consecutive elements of G'_n . Then there is an element of K_n in $(r/s, r'/s')$ having form 1 at r'/s' .

Minimal quadratics. In order to determine the least element of K_n greater than r/s (for fixed n) we need to introduce the concept of a minimal quadratic at r/s .

Let r/s be a non-negative element of G'_n ($r \geq 0, s > 0$), and let Δ be a positive integer. By (4) above, there is a solution $f = f_1$ in Z^3 of

$$(18) \quad f = (a, b, c), \quad f(r, s) = \Delta,$$

and f satisfies (18) if and only if

$$(19) \quad f - f_1 = \alpha(s, -r, 0) + \beta(0, s, -r),$$

where $\alpha, \beta \in Z$. (Of course (19) is valid for arbitrary r/s in G'_n .) From this, and the fact that $\Delta > 0$, it is easy to check that (18) has a solution f such that

$$(20) \quad a \geq -n, \quad |b| \leq n, \quad |c| \leq n,$$

(even with $a > 0$). If we begin with such an f and subtract $(s, -r, 0)$ or $(0, s, -r)$ repeatedly, subject to the condition that the f obtained always satisfies (20) then the

process eventually stops, and then at least one of the following holds:

$$(21) \quad \begin{cases} (i) & a < -n + s, & b < -n + s \\ (ii) & a < -n + s, & c > n - r \\ (iii) & b > n - r, & c > n - r. \end{cases}$$

Definition 6F.1. Let $r/s \geq 0$, $r/s \in G'_n$, and let Δ be a positive integer. Then the polynomial f satisfying (18), (20), and at least one part of (21) will be called the minimal quadratic of form Δ (at r/s , for given n), will be denoted by g_Δ , and will be said to be of type (i), (ii), or (iii) if (21)(i), (ii), or (iii) holds respectively.

Remarks. Despite its name, a minimal quadratic is not assumed to have degree 2. For any g_Δ of type (i) or type (ii), $|g_\Delta| \leq n$. However, a g_Δ of type (ii) may have $|g_\Delta| > n$, that is, $a > n$, though not for small Δ . For example, if $r/s \neq 0, 1$, then $r^2 + rs + s^2 \geq 7$, so if $\Delta \leq 6$, then $a \leq 0$, so certainly $|g_\Delta| \leq n$. The following lemma is also easily proved from the inequalities in (21).

Lemma 6F.3. Let $r/s \geq 0$, $r/s \in G'_n$, and let $g_1 = (a^{(1)}, b^{(1)}, c^{(1)})$ be minimal of form 1. Then $a^{(1)} < 0$ and $c^{(1)} > 0$.

We have shown that g_Δ exists for each positive integer Δ . Also, g_Δ is unique. This is immediate from the first part of the following lemma.

Lemma 6F.4. Let $r/s \geq 0$, $r/s \in G'_n$, and let Δ be a positive integer.

- (i) If f is any integral polynomial satisfying (18),
(20), then

$$(22) \quad f = g_{\Delta} + \alpha(s, -r, 0) + \beta(0, s, -r),$$

where α, β are non-negative integers.

(ii) If, in addition, $|f| \leq n$, then $|g_{\Delta}| \leq n$.

Proof. (i) is easy to check, and (ii) is immediate from (i).

Theorem 6F.5. Let $r/s, r'/s'$ be consecutive elements of G'_n with $0 \leq r/s < r'/s'$, and let Δ be a positive integer. Then

- (i) $|g_1| \leq n$, and g_1 has a zero θ_1 in K_n , and in $(r/s, r'/s')$,
- (ii) $g_{\Delta} = \Delta g_1 + \alpha(s, -r, 0) + \beta(0, s, -r)$, where α, β are non-negative integers,
- (iii) if there is an f satisfying the conditions of Lemma 6F.4, and having a zero θ in $(r/s, r'/s')$, then g_{Δ} has a zero θ_{Δ} in $(r/s, r'/s')$, and

$$\theta_1 \leq \theta_{\Delta} \leq \theta,$$

with equality only when the quadratics are equal up to a constant factor.

Proof. Part (i) follows from Lemma 6F.1, and Lemma 6F.4. (Actually, the fact that $|g_1| \leq n$ is also easily seen by a direct argument similar to that in the Remarks above.)

(ii) follows from the inequalities (21).

(iii) $\theta_{\Delta} \leq \theta$ follows from Lemma 6F.4 for if α, β are not both zero in (21), then the graph of g_{Δ} lies below that of f for $x > r/s$. Similarly, $\theta_1 \leq \theta_{\Delta}$ follows from (ii).

The following corollary is immediate:

Corollary 6F.6. Let $r/s, r'/s'$ be consecutive elements of G'_n . Then the least element of K_n in $(r/s, r'/s')$ has form 1 at r/s , and the greatest element of K_n in $(r/s, r'/s')$ has form 1 at r'/s' .

This result suggests an easy way of eliminating the most apparent discrepancy (mentioned in §6A) in the determinant corresponding to $\theta_1, \theta_2, \theta_3$ when θ_2 is rational. We simply eliminate these determinants by inserting an extra (reducible) quadratic $(s, -r, 0)$ at the rational point r/s just before $(0, s, -r)$. To say that g_1 has form 1 at r/s is to say that on writing

$$g_1 = (a^{(1)}, b^{(1)}, c^{(1)}),$$

we have

$$\begin{vmatrix} s & -r & 0 \\ 0 & s & -r \\ a^{(1)} & b^{(1)} & c^{(1)} \end{vmatrix} = 1.$$

This is one of the two new determinants created by inserting $(s, -r, 0)$. In the case where $r/s \neq 0$, if (c', b', a') is the primitive minimal polynomial of form 1 at s/r of the least element of K_n greater than s/r , then (a', b', c') is the primitive minimal polynomial of form 1 at r/s of the greatest element of K_n less than r/s , and the other new determinant created is

$$\begin{vmatrix} a' & b' & c' \\ s & -r & 0 \\ 0 & s & -r \end{vmatrix} = \begin{vmatrix} r & -s & 0 \\ 0 & r & -s \\ c' & b' & a' \end{vmatrix} = 1.$$

(The case $r/s = 0$ can be handled separately, the least positive element of K_n being given by $(a^{(1)}, b^{(1)}, c^{(1)}) =$

$(-n, -n, 1)$ and the greatest negative element being given by $(a', b', c') = (-n, n, 1)$.)

The eliminated determinant was

$$D = \begin{vmatrix} a' & b' & c' \\ 0 & s & -r \\ a^{(1)} & b^{(1)} & c^{(1)} \end{vmatrix}.$$

Now by Lemma 6F.4(i),

$$(23) \quad (a', b', c') = (a^{(1)}, b^{(1)}, c^{(1)}) + \alpha(s, -r, 0) + \beta(0, s, -r),$$

for non-negative integers α, β , not both zero. Hence

$D = \alpha$. Also, when $r/s > 0$, it follows from Lemma 6F.3 applied to g_1 , and to (c', b', a') , that $a^{(1)} < 0$ and $a' > 0$. Then from (23), D is positive, as was observed in the numerical data of Brown and Mahler [2] in the reference to Table 2. (Their determinants were actually negative, as their quadratics had non-negative leading coefficients.)

§6G. The second smallest element of K_n in $(r/s, r'/s')$.

In this section, we will suppose that $r/s, r'/s'$ are consecutive elements of G'_n with $0 \leq r/s < r'/s'$, although the results may be easily translated into results for negative r/s . We will let

$$f^{(i)} = (a^{(i)}, b^{(i)}, c^{(i)}) \quad (i = 1, 2)$$

denote the primitive quadratics of the smallest and second smallest elements $\theta^{(1)}, \theta^{(2)}$ in $(r/s, r'/s')$. Thus, Corollary 6F.6 shows that $f^{(1)} = g_1$. I have also investigated the quadratic $f^{(2)}$ by methods which include further study of the minimal quadratics of §6F. As the details are lengthy I will not include them here, but

I shall describe the results and show how they can be used in connection with the determinant

$$\det((0, s, -r), f^{(1)}, f^{(2)}) .$$

It will be convenient to express an arbitrary integral quadratic f such that

$$f(r, s) = \Delta > 0$$

in the notation

$$\begin{aligned} f &= [\Delta, \alpha, \beta] \\ &= \Delta f^{(1)} + \alpha(s, -r, 0) + \beta(0, s, -r) \end{aligned}$$

where by Lemma 6F.4(i), α, β are non-negative integers if $|f| \leq n$. This notation is useful not only in determining, of two given quadratics, which has the least zero greater than r/s , but also in applying various techniques of obtaining from a given quadratic, a quadratic with a lesser zero in $(r/s, r'/s')$.

For example, if

$$f = [\Delta, \alpha, \beta] ,$$

where $\beta = 0$, $\alpha \geq 2$, is the primitive minimal polynomial of an element θ_f of K_n in $(r/s, r'/s')$, and $[x]^+$ denotes the least integer not less than x , and

$$h = \left[\left[\frac{\Delta}{\alpha} \right]^+, 1, 0 \right] ,$$

then $|h| \leq n$, and h has a zero θ_h in $(r/s, r'/s')$, such that

$$\theta^{(1)} < \theta_h < \theta_f ,$$

so $\theta_f \neq \theta^{(2)}$, and $f \neq f^{(2)}$.

Let us write in square bracket notation,

$$f^{(2)} = [\Delta^{(2)}, \alpha^{(2)}, \beta^{(2)}] .$$

Then it is easily checked that

$$(1F) \quad D = \det((0, s, -r), f^{(1)}, f^{(2)}) = \alpha^{(2)} .$$

I have proved by detailed considerations, that

$$(2) \quad f^{(2)} = [\Delta, 1, 0], [\Delta, 0, 1], [\Delta, 1, 1], \\ [3, 1, 2], \text{ or } [3, 2, 1],$$

and that the case $f^{(2)} = [3, 2, 1]$ can only occur when $r/s > 1$ and $b^{(1)} + s > n$. It follows that the determinant in (1) satisfies

$$(3) \quad D = 0, 1, 2, \text{ and } D = 0, 1 \text{ if } 0 \leq r/s \leq 1.$$

(Negative determinants can occur if we normalise the quadratics to get non-negative leading coefficients.)

The possibility $f^{(2)} = [3, 1, 2]$ is interesting. An example where this holds is when $n = 11$, $r/s = 8/11$, $r'/s' = 3/4$, when it can be checked by considering minimal quadratics in the sense of §6F and using the various reduction techniques alluded to in the present section, that

$$f^{(1)} = [1, 0, 0] = (-6, -8, 9), \\ f^{(2)} = [3, 1, 2] = (-7, -10, 11).$$

A consequence of this example is that at $11/8$ in K_{11} we have the two consecutive quadratics preceding $11/8$, forming with $11/8$, the determinant,

$$\begin{vmatrix} 11 & -10 & -7 \\ 9 & -8 & -6 \\ 0 & 8 & -11 \end{vmatrix} = -2 \quad (= -\beta^{(2)} \text{ at } 8/11)$$

Thus the determinant property for three consecutive elements of K_n does not hold in this example where the largest element is rational (providing another counter-example to a conjecture of Brown and Mahler), although the determinant property does hold (at this point) in the modified scheme (explained after Corollary 6F.5) where we place $(8, -11, 0)$ before $(0, 8, -11)$, for then the determinant becomes

$$\begin{vmatrix} 11 & -10 & -7 \\ 9 & -8 & -6 \\ 8 & -11 & 0 \end{vmatrix} = -1 \quad (= -\alpha^{(2)} \text{ at } 8/11).$$

I have not determined whether the case $f^{(2)} = [3, 2, 1]$ can actually occur (when $r/s > 1$). If it cannot, then it would follow that the determinant property for three consecutive elements

$$f, \quad g, \quad (s, -r, 0)$$

in our modified scheme also holds (for arbitrary $r/s \in G'_n$).

§6H. The quadratics of form 1 at r/s .

In this section I will prove the determinant property for quadratics of form 1. That is, I will prove the following theorem.

Theorem 6H.1. Let $r/s, r'/s'$ be consecutive elements of G'_n . Then the primitive quadratics of elements of K_n in $(r/s, r'/s')$ of form 1 have the property that the determinant of any 3 consecutive quadratics is $0, \pm 1$.

Proof. Let S be the set of all integral quadratics

$$f(x) = ax^2 + bx + c$$

such that

$$f(r, s) = 1, \quad |f| \leq n,$$

and let S_1 be the subset of S consisting of the quadratics in S satisfying

$$f(r', s') < 0.$$

By Lemmas 6B.1(ii) and 6B.2(i), there is a one-one correspondence between the elements of K_n in $(r/s, r'/s')$ and the quadratics in S_1 . We will prove the determinant result about S_1 by proving a determinant result about S .

As we found in §6F(19) (which is of course valid for arbitrary r/s in G'_n), the elements of S may be written as

$$(1) \quad f = f_0 + \alpha(s, -r, 0) + \beta(0, s, -r), \quad \alpha, \beta \in \mathbb{Z},$$

where f_0 is a particular quadratic of form 1 at r/s .

For $i = 1, 2$ and elements f_i of S such that

$$(2) \quad f_i = f_0 + \alpha_i(s, -r, 0) + \beta_i(0, s, -r)$$

we have

$$(3) \quad f_1(x) - f_2(x) = (sx - r) \{ (\alpha_1 - \alpha_2)x + (\beta_1 - \beta_2) \}.$$

Thus the graphs of $f_1(x)$, $f_2(x)$ cross at r/s , and if $\alpha_1 \neq \alpha_2$, at $-(\beta_1 - \beta_2)/(\alpha_1 - \alpha_2)$. Thus the graphs cross in $(r/s, r'/s')$ if and only if

$$(4) \quad -\frac{\beta_1 - \beta_2}{\alpha_1 - \alpha_2} \in (r/s, r'/s').$$

We now show that this is not possible.

Without loss of generality, suppose that $r/s \geq 0$, $s > 0$. Since $r/s, r'/s'$ are consecutive in K_n , the larger of $r+r'$, $s+s'$ exceeds n . Suppose that (4) holds. Then

$$|\beta_1 - \beta_2| \geq r+r', \quad |\alpha_1 - \alpha_2| \geq s+s'.$$

If $r \geq s$, then $r' > s'$, $r+r' > s+s'$, so $r+r' > n$, $|\beta_1 - \beta_2| > n$. By considering the third components of f_1, f_2 , we see that $r \leq 1$, $s = 1$, $r/s = 1/1$. On the other hand, if $r < s$, then $r' \leq s'$, $r+r' < s+s'$, $s+s' > n$, $|\alpha_1 - \alpha_2| > n$, and from the first components of f_1, f_2 , $s \leq 1$, so $r/s = 0/1$. Thus $(r/s, r'/s') = \left(\frac{0}{1}, \frac{1}{n}\right)$ or $\left(\frac{1}{1}, \frac{n}{n-1}\right)$. But it is readily seen from §6D, §6E respectively that (4) is false in these cases. (In fact the determinant problem has been disposed of for these intervals.) Thus we have shown that (4) does not hold, that is the graphs of the

quadratics f in S do not cross in $(r/s, r'/s')$, and so an ordering on the quadratics in S may be defined as follows: If $f_1, f_2 \in S$, then

$$(5) \quad f_1 < f_2 \Leftrightarrow f_1(x) < f_2(x) \quad \forall x \in (r/s, r'/s').$$

Clearly, this ordering on S agrees with the ordering on S_1 corresponding to the order of the zeros of the quadratics in S_1 , and S_1 is an initial segment of S in the sense that there is an element f_N in S_1 such that

$$S_1 = \{f \in S \mid f \leq f_N\}.$$

Thus, we need only prove the corresponding determinant result for S .

It is also clear from (3) that if f_1, f_2 in S are given by (2), then

$$(6) \quad f_1 < f_2 \Leftrightarrow \alpha_1 r + \beta_1 s < \alpha_2 r + \beta_2 s$$

$$\text{or } \begin{cases} \alpha_1 r + \beta_1 s = \alpha_2 r + \beta_2 s \\ \text{and } \alpha_1 < \alpha_2 \end{cases}.$$

Thus, the ordering $<$ on S is a "linear ordering". One final observation about S which will make our proof work is that S consists of the points (a, b, c) of the lattice \mathbb{Z}^3 which lie in a 2-dimensional convex set, namely that given by the inequalities

$$\begin{cases} -n \leq a = a_0 + \alpha s \leq n \\ -n \leq b = b_0 - \alpha r + \beta s \leq n \\ -n \leq c = c_0 - \beta r \leq n, \end{cases}$$

or, in terms of α, β , the point $\alpha = (\alpha, \beta)$ lies in the convex set K in \mathbb{R}^2 given by the inequalities

$$\begin{cases} (-n - a_0)/s \leq \alpha \leq (n - a_0)/s \\ -n - b_0 \leq -\alpha r + \beta s \leq n - b_0 \\ (-n - c_0)/r \leq \beta \leq (n - c_0)/r. \end{cases}$$

Let f_1, f_2, f_3 be consecutive elements of S given by (2) with $i = 1, 2, 3$ and write also

$$(7) \quad f_i = (a_i, b_i, c_i) \quad (i = 1, 2, 3).$$

Then the relevant determinant is

$$\begin{aligned} D &= \det \begin{bmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{bmatrix} \\ &= \det \left(\begin{bmatrix} 1 & \alpha_1 & \beta_1 \\ 1 & \alpha_2 & \beta_2 \\ 1 & \alpha_3 & \beta_3 \end{bmatrix} \begin{bmatrix} a_0 & b_0 & c_0 \\ s & -r & 0 \\ 0 & s & -r \end{bmatrix} \right) \\ &= \det \begin{bmatrix} 1 & \alpha_1 & \beta_1 \\ 1 & \alpha_2 & \beta_2 \\ 1 & \alpha_3 & \beta_3 \end{bmatrix} \\ &= \det \begin{bmatrix} \alpha_2 - \alpha_1 & \beta_2 - \beta_1 \\ \alpha_3 - \alpha_2 & \beta_3 - \beta_2 \end{bmatrix} \\ &= \det (\alpha_2 - \alpha_1, \alpha_3 - \alpha_2), \end{aligned}$$

where $\alpha_i = (\alpha_i, \beta_i)$ for $i = 1, 2, 3$.

Suppose that $|D| > 1$. Then by Lemma 3B.5, or by Theorem 34 of Hardy and Wright [6], there is a point $\alpha_4 = (\alpha_4, \beta_4)$ in Z^2 lying within the triangle with vertices $\alpha_1, \alpha_2, \alpha_3$ or on its boundary, but distinct from α_i for $i \neq 4$. Since $\alpha_4 \in K$, it follows from (6) that the quadratic f_4 of S , given by

$$f_4 = f_0 + \alpha_4(s, -r, 0) + \beta_4(0, s, -r),$$

lies between f_1 and f_3 . But $f_4 \neq f_2$, and this contradicts the consecutiveness of f_1, f_2, f_3 . Thus $|D| \leq 1$, as required. This completes the proof of Theorem 6H.1.

CHAPTER 7

AN INDEX PROBLEM OF CASSELS§7A Introduction.

Throughout this chapter, M will denote a 3-dimensional lattice, and L a 2-dimensional lattice in R^3 , and the octants in R^3 will be denoted in the obvious way by A_{+++} , A_{-++} , and so on. For example

$$A_{-++} = \{\underline{x} \mid x_1 \leq 0, x_2 \geq 0, x_3 \geq 0\}.$$

Similarly the quadrants in R^2 will be denoted by $\pm A_{++}$, $\pm A_{-+}$.

Throughout this chapter, we will also let F denote the convex distance function in R^3 given by

$$(1) \quad F(\underline{x}) = F(x_1, x_2, x_3) = |x_1| + |x_2| + |x_3|.$$

Now let M be a 3-dimensional lattice, and A an octant in R^3 . Define a minimal point of M in A with respect to F as a point \underline{u} in $\dot{M} \cap A$ with minimal $F(\underline{u})$.

For simplicity, suppose that \dot{M} contains no points on the coordinate planes, and let X be a set of 4 minimal points, one for each pair of opposite octants $\pm A$. Then we have the following problem, which is the subject of the present chapter.

Index problem: What can be said about the index $\{M: \ell(X)\}$ of the sublattice of M generated by X ?

This problem was raised by Professor J.W.S. Cassels in the Arithmetic Seminar in Cambridge in 1974. Professor Cassels suggested that one ought to be able to show that the index is at most 6, or perhaps at most 2, or perhaps even always equal to 1. The existence of a small bound for the index is suggested by results on the classical successive minima. For

example, by VIII.4.3 Corollary of Cassels [3], the index of the associated minimal vectors in R^n does not exceed $n!$ In fact, for $n = 3$, Bantegnie [1] has shown that the index of these minimal vectors does not exceed 4. (See also Lekkerkerker [8] p.56.)

For the problem as stated above, the index need in fact not be 1. In §7B we will give an example where the index is 2, and to retain the possibility that the index be always equal to 1, we will introduce one modified 'minimal' lattice point. We will also discuss briefly the connection with successive minima.

Then in §7C, I will prove index results for 2-dimensional lattices in R^2 or R^3 . Sections §7D, §7E concern the two main ideas which we will use to tackle the index problem of Cassels. Finally, in §7F I will obtain a bound for this index.

§7B. The set of minimal points of a 3-dimensional lattice, notation, examples.

Let M be a 3-dimensional lattice in R^3 . Suppose that \dot{M} has no points on the coordinate planes. Then each point \underline{u} in \dot{M} belongs to a uniquely determined octant $A(\underline{u})$. Let F be as in §7A(1).

Let the term minimal point of M in A (or of $\dot{M} \cap A$) be as defined in §7A, and let X be a set of minimal points as in §7A. Suppose that

$$(1) \quad X = \{\underline{a}, \underline{b}, \underline{c}, \underline{d}\},$$

where

$$(2) \quad F(\underline{a}) \leq F(\underline{b}) \leq F(\underline{c}), \text{ and } F(\underline{b}) \leq F(\underline{d}).$$

Since \dot{M} has no points on the coordinate planes, $\underline{a}, \underline{b}$ are

linearly independent. Let $\underline{c}', \underline{d}'$ be points such that

$$(3) \quad \begin{cases} \underline{c}' \in M \cap A(\underline{c}) \sim \text{lin}(\underline{a}, \underline{b}) = V \text{ say,} \\ \underline{d}' \in M \cap A(\underline{d}) \sim \text{lin}(\underline{a}, \underline{b}) = W \text{ say,} \\ F(\underline{c}') = \min \{F(\underline{u}) \mid \underline{u} \in V\}, \\ F(\underline{d}') = \min \{F(\underline{u}) \mid \underline{u} \in W\}, \end{cases}$$

and let

$$X' = \{\underline{a}, \underline{b}, \underline{c}', \underline{d}'\}.$$

In this chapter, we shall investigate both indices

$$(M:l(X)), \quad (M:l(X')).$$

Since a plane in R^3 through $\underline{0}$ contains interior points of at most three octant pairs $\pm A$, it follows that at most one of \underline{c} and \underline{d} lies in $\text{lin}(\underline{a}, \underline{b})$. Hence we may suppose that $\underline{c} \notin \text{lin}(\underline{a}, \underline{b})$, so that

$$(4) \quad \begin{aligned} \underline{c}' &= \underline{c}, \quad X' = \{\underline{a}, \underline{b}, \underline{c}, \underline{d}'\}, \\ X' &= X \Leftrightarrow \underline{d} \in \text{lin}(\underline{a}, \underline{b}). \end{aligned}$$

We now give an example where Cassels's index $(M:l(X))$ is 2 but the modified index $(M:l(X'))$ is 1.

Example 1.

We note first that if \underline{x} belongs to the lattice with basis

$$(0, 1, 0), \quad (0, 0, 1), \quad (1, \tfrac{1}{2}, \tfrac{1}{2})$$

and \underline{x} is not a linear combination of $(0, 1, 0)$ and $(0, 0, 1)$, then $F(\underline{x}) \geq 2$. We consider a lattice M which is a slight deformation of the above lattice, namely the lattice M with basis

$$\underline{u} = (-\varepsilon_1, 1, \varepsilon_2),$$

$$\underline{v} = (-\varepsilon_3, -\varepsilon_4, 1),$$

$$\underline{w} = (1, \tfrac{1}{2} + \varepsilon_5, \tfrac{1}{2} + \varepsilon_6)$$

where the ε_i are small positive numbers chosen to satisfy

conditions to be stated. The ε_i are to satisfy the appropriate linear independence conditions over the rationals, so that \dot{M} contains no point on a coordinate plane. The ε_i are to be small enough so that $F(\underline{x}) > 1\frac{1}{2}$ say for all \underline{x} in $M \sim \ell(\underline{u}, \underline{v})$. It then follows that one of the minimal points $\underline{a}, \underline{b}$ (see (2)) is \underline{u} in A_{-++} , and the other is \underline{v} in A_{--+} , and that these are unique. (All the while we are assuming that the ε_i are small enough.)

Now let $\underline{c}, \underline{d}$ be minimal points of M in A_{+++} and A_{+--} respectively. Consider the vector

$$\underline{u} \times \underline{v} = (1 + \varepsilon_2 \varepsilon_4, -\varepsilon_2 \varepsilon_3 + \varepsilon_1, \varepsilon_1 \varepsilon_4 + \varepsilon_3),$$

normal to $\text{lin}(\underline{u}, \underline{v})$, and consider the lattice point

$$\underline{v} - \underline{u} = (-\varepsilon_3 + \varepsilon_1, -\varepsilon_4 - 1, 1 - \varepsilon_2).$$

Suppose that

$$(5) \quad \varepsilon_1 > \varepsilon_3, \quad \varepsilon_2 < 1.$$

Then

$$\underline{u} \times \underline{v} \in A_{+++}, \quad \underline{v} - \underline{u} \in A_{+--},$$

and $\text{lin}(\underline{u}, \underline{v})$ contains no point of $\dot{M} \cap A_{+++}$. Hence

$$\underline{c} \notin \text{lin}(\underline{u}, \underline{v}) = \text{lin}(\underline{a}, \underline{b}),$$

$$\underline{c}' = \underline{c},$$

where \underline{c}' is as in (3).

Let us consider the minimal values of $F(\underline{x})$, where \underline{x} lies in a coset

$$(6) \quad \begin{cases} \underline{x} \in h\underline{w} + \ell(\underline{u}, \underline{v}) & (h \in \mathbb{Z}), \\ \underline{x} \neq \pm \underline{u}, \pm \underline{v}, 0 \\ \underline{x} \in A_{+++} \cup A_{+--}. \end{cases}$$

For $h = 0$, we have

$$(7) \quad \underline{v} - \underline{u} \in A_{+--}, \quad F(\underline{v} - \underline{u}) = 2 + \varepsilon_1 - \varepsilon_3 + \varepsilon_4 - \varepsilon_2,$$

and

$$F(\underline{x}) > 2\frac{1}{2} \text{ for other } \underline{x} \text{ with } h = 0.$$

For $h = 1$, we have

$$(8) \quad \underline{w} \in A_{+++}, \quad F(\underline{w}) = 2 + \varepsilon_5 + \varepsilon_6.$$

We also have

$$(9) \quad \begin{cases} \underline{w} - \underline{u} = (1 + \varepsilon_1, -\frac{1}{2} + \varepsilon_5, \frac{1}{2} + \varepsilon_6 - \varepsilon_2) \\ \underline{w} - \underline{u} \in A_{+-+}, \\ F(\underline{w} - \underline{u}) = 2 + \varepsilon_1 - \varepsilon_5 + \varepsilon_6 - \varepsilon_2, \end{cases}$$

and

$$F(\underline{x}) > 2\frac{1}{2} \text{ for other } \underline{x} \text{ with } h = 1.$$

For $h = 2$, we have

$$2\underline{w} - \underline{u} - \underline{v} = (2 + \varepsilon_1 + \varepsilon_3, 2\varepsilon_5 + \varepsilon_4, 2\varepsilon_6 - \varepsilon_2)$$

which lies in A_{+++} on assuming that

$$(10) \quad 2\varepsilon_6 > \varepsilon_2,$$

and we have

$$(11) \quad F(2\underline{w} - \underline{u} - \underline{v}) = 2 + \varepsilon_1 + \varepsilon_3 + 2\varepsilon_5 + \varepsilon_4 + (2\varepsilon_6 - \varepsilon_2)$$

and

$$F(\underline{x}) > 2\frac{1}{2} \text{ for other } \underline{x} \text{ with } h = 2.$$

For $h \geq 3$, we have $F(\underline{x}) > 3\frac{1}{2}$, and for $h \leq 0$, \underline{x} clearly cannot lie in $A_{+++} \cup A_{+-+}$ unless $F(\underline{x})$ is rather large, certainly greater than $3\frac{1}{2}$ (always assuming that the ε_i are small enough).

Thus, in considering the conditions (6), we find that the minimal point $c = \underline{c}'$ in A_{+++} must be one of the points \underline{w} , $2\underline{w} - \underline{u} - \underline{v}$ which satisfy (8), (10) and (11), and which have $h = 1, 2$ respectively. We can ensure that

$$(12) \quad \underline{c} = 2\underline{w} - \underline{u} - \underline{v}$$

by choosing the ε_i so that

$$F(2\underline{w} - \underline{u} - \underline{v}) < F(\underline{w}),$$

that is, to satisfy

$$(13) \quad \varepsilon_1 + \varepsilon_3 + \varepsilon_5 + \varepsilon_4 + \varepsilon_6 > \varepsilon_2 ,$$

and (13) is certainly compatible with (5) and (10).

Equation (12) implies that

$$(14) \quad (M: \ell(\underline{a}, \underline{b}, \underline{c})) = 2.$$

From the above consideration of the conditions (6), we also conclude that the point \underline{d}' in A_{+-+} (see (3)) is

$$(15) \quad \underline{d}' = \underline{w} - \underline{u}$$

calculated in (9), so that with X' as in (4),

$$(M: \ell(X')) = 1 ,$$

as claimed.

We next prove that the minimal point \underline{d} in A_{+-+} is

$$(16) \quad \underline{d} = \underline{v} - \underline{u} ,$$

(see (7)).

Inequalities (10), (13) imply that

$$\begin{aligned} \varepsilon_6 &> \varepsilon_1 + \varepsilon_3 + \varepsilon_5 + \varepsilon_4 \\ &> \varepsilon_5 + \varepsilon_4 - \varepsilon_3 . \end{aligned}$$

Hence by (7), (9),

$$F(\underline{v} - \underline{u}) < F(\underline{w} - \underline{u}) ,$$

and (16) follows. From this, we see that with X as in (1),

$$\begin{aligned} \ell(X) &= \ell(\underline{a}, \underline{b}, \underline{c}, \underline{d}) \\ &= \ell(\underline{a}, \underline{b}, \underline{c}) , \end{aligned}$$

so by (14), we have

$$(M: \ell(X)) = 2 ,$$

as we claimed.

Connection with successive minima.

Let $\underline{a}, \underline{b}, \underline{c}, \underline{d}$ be as in (1), (2). It is clear that the first minimal point \underline{a} is just a first successive minimal point of the lattice M with respect to the distance function

F. (See Cassels [3], p.201.) Further, given a lattice M , it may well turn out that the successive minimal points are $\underline{a}, \underline{b}, \underline{c}$ or $\underline{a}, \underline{b}, \underline{d}$. However, we now give a simple example which illustrates that the second successive minimal point need not even lie in $\text{lin}(\underline{a}, \underline{b})$.

Example 2.

Let M be the lattice in R^3 with basis

$$u = (1, \varepsilon_1, \varepsilon_2),$$

$$v = (\varepsilon_3, 1, \varepsilon_4),$$

$$w = (\varepsilon_5, \varepsilon_6, -1\frac{1}{2}),$$

where

$$\varepsilon_1 + \varepsilon_2 < \varepsilon_3 + \varepsilon_4,$$

$$0 < \varepsilon_i < \frac{1}{8} \quad \text{for all } i.$$

Then $\underline{u}, \underline{v}, \underline{w}$ are successive minimal points, and $\underline{a} = \underline{u}$, but $\underline{b} = \underline{w}$, so

$$\text{lin}(\underline{a}, \underline{b}) \neq \text{lin}(\underline{u}, \underline{v}).$$

Let M be a 3-dimensional lattice in R^3 , let λ_1, λ_2 denote the first two successive minima of M with respect to F , and let $\underline{a}, \underline{b}$ denote the first two minimal points as in (1), (2). Also, let

$$(17) \quad F(\underline{a}) = \mu_1, \quad F(\underline{b}) = \mu_2.$$

Then

$$(18) \quad \lambda_1 = \mu_1, \quad \mu_2 - \lambda_1 \leq \lambda_2 \leq \mu_2.$$

This is trivial, except for the left hand inequality, which will be proved in §7D.

In the case when the (classical) second successive minimal point \underline{m}_2 of M lies in $\text{lin}(\underline{a}, \underline{b})$ it may be proved that, re-choosing \underline{b} if necessary,

$$(19) \quad \underline{b} = \underline{m}_2 \quad \text{or} \quad \underline{m}_2 - \underline{a}$$

(the proof would use Theorems 7C.1, 7C.3.), and in this case,

(18) obviously follows from (19).

§7C. Minimal points of a 2-dimensional lattice.

In this section, we let L be a general 2-dimensional lattice in R^3 which may have points on the coordinate planes or axes. We will prove if $\underline{a}, \underline{b}$ are the minimal points as defined below, then

$$\{L : \lambda(\underline{a}, \underline{b})\} = 1,$$

except in certain cases when L contains points on the coordinate axes, when the index is 2. It will then follow that the minimal points $\underline{a}, \underline{b}$ of §7B (1), (2) are a basis for $M \cap \text{lin}(\underline{a}, \underline{b})$.

Let P be an arbitrary plane in R^3 passing through the origin. There are three possibilities:

- (i) P is a coordinate plane,
- (ii) P is not a coordinate plane but passes through a coordinate axis.
- (iii) P contains no coordinate axis.

These cases correspond to the cases when the Cartesian equation for P has two, one or none of its coefficients equal to zero. Consider the intersections $P \cap A$, where A is an octant in R^3 . When non-degenerate, $P \cap A$ is a 2-dimensional cone, which we will refer to as a coordinate sector in P . In cases (i), (ii) P has four distinct sectors. (In case (ii), exactly two of the lines of intersection of P with a coordinate plane coincide.) In case (iii), P has six distinct sectors, these being determined by the three distinct lines in which the coordinate planes intersect P .

Now consider the convex distance function F of §7A (1). That is,

$$(1) \quad F(\underline{x}) = |x_1| + |x_2| + |x_3|.$$

We will be using the fact that F is linear in each octant and hence in each sector in P . We will also use the fact that the inequality

$$(2) \quad F(\tilde{x} + \tilde{y}) < F(\tilde{x}) + F(\tilde{y})$$

holds whenever no octant contains both \tilde{x} and \tilde{y} , and so holds for any two points of P which do not lie in a common sector.

The minimal points \tilde{a}, \tilde{b} of L .

Let L be an arbitrary 2-dimensional lattice in the plane P in R^3 . Choose \tilde{a} in L to satisfy

$$(3) \quad \tilde{a} \in \dot{L}, \quad F(\tilde{a}) = \min_{\tilde{u} \in \dot{L}} F(\tilde{u}).$$

Choose a sector A in P such that $\tilde{a} \in A$, and choose \tilde{b} in L such that

$$(4) \quad \tilde{b} \in L \sim (\text{relint } A) \sim (-\text{relint } A) \sim \ell(\tilde{a}) = T, \quad \text{say,}$$

$$(5) \quad F(\tilde{b}) = \min_{\tilde{u} \in T} F(\tilde{u}).$$

We will refer to \tilde{a}, \tilde{b} as the first and second minimal points of L with respect to F . (An analogous general definition could be given for the first two minimal points of an arbitrary 3-dimensional lattice, though we have omitted this for the sake of simplicity.) If M is a 3-dimensional lattice in R^3 such that \dot{M} has no points on the coordinate planes, if \tilde{a}, \tilde{b} are taken as the first two minimal points as in §7B, and if $L = M \cap \text{lin}(\tilde{a}, \tilde{b})$, then \tilde{a}, \tilde{b} are minimal points of L in the above sense, as is easily seen.

Theorem 7C.1. Let L be a 2-dimensional lattice in the plane P in R^3 , and let \tilde{a}, \tilde{b} be minimal points with respect to F as given by (1), (3), (4), (5) above. Then

$$(L : \ell(\tilde{a}, \tilde{b})) = 1$$

except when $F(\underline{a}) = F(\underline{b})$ and L has basis $\underline{a}, \frac{1}{2}(\underline{a}+\underline{b})$ (in which case the index is 2). The exceptional case can only occur when P has just four sectors, one of which contains $\underline{a}, \underline{b}$ on its bounding rays.

Note 1: Thus, in the exceptional case, one of $\underline{a}, \underline{b}$ lies on a coordinate axis, and the other lies on a coordinate plane.

Note 2: If $P = \{\underline{x} | x_3 = 0\}$ we may identify F with the corresponding distance function on R^2 , and infer the analogous result for 2-dimensional lattices in R^2 .

Proof of Theorem. Let F_A denote the linear function that F coincides with on A . Replacing \underline{b} by $-\underline{b}$ if necessary we may suppose that the following (strict) inequality holds

$$(6) \quad F_A(\underline{b}) < F(\underline{a})$$

(that is, $\underline{0}, \underline{b}$ lie on the same side of the line $F_A(\underline{x}) = F(\underline{a})$ in P). From (3), we have

$$(7) \quad F(\underline{b}) \geq F(\underline{a}).$$

We suppose the required index is not 1. Then by Lemma 3B.5 the triangle $\underline{0} \underline{a} \underline{b}$ contains a point \underline{y} of L other than $\underline{0}, \underline{a}, \underline{b}$; any such point is of the form

$$\underline{y} = \lambda \underline{a} + \mu \underline{b}, \quad 0 \leq \lambda < 1, \quad 0 \leq \mu < 1, \quad 0 < \lambda + \mu \leq 1.$$

Now since $\underline{b} \notin \text{relint } A$, $\underline{a} \in A$, we have that $\underline{y} \notin \text{relint } A$. Thus it follows from (3), (4), (6) and the linearity of F_A that \underline{y} must be in T as defined by (4). If $F(\underline{a}) < F(\underline{b})$, or \underline{a} and \underline{b} are not contained in a common sector of P , or $\lambda + \mu < 1$, then by (2), (7) and the convexity of F we obtain $F(\underline{y}) < F(\underline{b})$, a contradiction to (5). Hence we must have

- (i) $F(\underline{b}) = F(\underline{a})$,
- (ii) some sector in P contains both \underline{a} and \underline{b} , and
- (iii) \underline{y} is on the line-segment $\underline{a}\underline{b}$.

It follows from (iii) that $\underline{b} - \underline{a}$ is not primitive, so let h be the largest integer such that

$$(8) \quad \underline{a} - \underline{b} = h\underline{u} \quad (h \geq 2)$$

for some \underline{u} in L . By (3), we have

$$F(\underline{a}-\underline{b}) \geq hF(\underline{a}).$$

But

$$\begin{aligned} F(\underline{a}-\underline{b}) &\leq F(\underline{a}) + F(-\underline{b}) \\ &= 2F(\underline{a}) \quad \text{by (i).} \end{aligned}$$

Hence

$$(9) \quad F(\underline{a}+(-\underline{b})) = F(\underline{a}) + F(-\underline{b}),$$

and

$$(10) \quad h = 2.$$

From (9) (10) we have that

- (iv) some sector in P contains both \underline{a} and $-\underline{b}$.

Now conditions (ii), (iv) imply that there are two adjacent sectors in P which both contain \underline{a} , and which each contain one of \underline{b} , $-\underline{b}$. This can only happen when P has only four sectors, and \underline{a} , \underline{b} lie on the bounding rays of one sector. Moreover, (iii) and (10) show that $\underline{y} = \frac{1}{2}(\underline{a}+\underline{b})$, and that \underline{a} , \underline{y} form a basis for L . Thus if the index is not 1 then the exceptional case specified in the theorem must occur.

We now have the following immediate corollary.

Theorem 7C.2. Let M be a 3-dimensional lattice in R^3 having no non-zero points on the coordinate planes, and let \underline{a} , \underline{b} , \underline{c} , \underline{d} , \underline{x} be as given in §7B (1), (2), (3), (4). Then

- (i) $M \cap \text{lin}(\underline{a}, \underline{b}) = \ell(\underline{a}, \underline{b}),$
(ii) If $\underline{c} \in \text{lin}(\underline{a}, \underline{b}),$ then $\ell(X) = \ell(\underline{a}, \underline{b}, \underline{c}).$

The following result can be proved by a proof similar to that of Theorem 7C.1, but simpler, and we omit the details.

Theorem 7C.3. Let M be a 3-dimensional lattice in R^3 with first two (classical) successive minimal points $\underline{a}, \underline{b}$ with respect to F . Then the index

$$(M \cap \text{lin}(\underline{a}, \underline{b}) : \ell(\underline{a}, \underline{b}))$$

is only greater than 1 in the same exceptional cases as in Theorem 7C.1 with $L = \text{lin}(\underline{a}, \underline{b}).$

In Theorem 7C.3, the exceptional cases of course cannot occur if \dot{M} has no points on the coordinate planes.

§7D. Distance of lattice planes from the origin.

In this section, we will let M be a 3-dimensional lattice in R^3 , and although not absolutely necessary, we will assume, without explicit mention, that \dot{M} has no points on the coordinate planes. We will consider a lattice plane not through the origin, and parallel to $\text{lin}(\underline{a}, \underline{b})$ where $\underline{a}, \underline{b}$ are as in §7B, and find a lower bound for its distance from $\underline{0}$ with respect to the distance function F of §7A(1).

Consider a plane P in R^3 , with equation

$$(1) \quad \gamma_1 x_1 + \gamma_2 x_2 + \gamma_3 x_3 = t, \quad |\gamma_1| = \max_i |\gamma_i|.$$

By the distance τ of the plane P from $\underline{0}$ with respect to F we will mean the minimum value of F on P . We have

$$(2) \quad \tau = F(t/\gamma_1, 0, 0) = |t/\gamma_1|.$$

This is easily seen from the fact that F is non-negative and linear on each octant A , and so achieves a minimum value at a vertex of one of the polygonal regions $P \cap A$.

The following lemma establishes (18) of §7B.

Lemma 7D.1. Let M be a 3-dimensional lattice in R^3 , let \underline{a} and \underline{b} be as in §7B(1), (2), $\mu_1 = F(\underline{a})$, $\mu_2 = F(\underline{b})$ and let $\lambda_1 (= \mu_1)$, λ_2 be the first two successive minima as in §7B(17), (18). Then

$$(3) \quad \lambda_2 \geq \max(\mu_2 - \mu_1, \mu_1).$$

Proof. The only non-trivial part to prove is

$$(4) \quad \lambda_2 \geq \mu_2 - \mu_1.$$

Let \underline{u} be a point of \dot{M} such that

$$(5) \quad F(\underline{u}) < \mu_2 - \mu_1.$$

We must prove that

$$(6) \quad \underline{u} \in \ell(\underline{a}).$$

Let A be the octant such that $\underline{a} \in A$.

By definition, we have

$$(7) \quad F(\underline{a}) = \mu_1 = \min_{\underline{v} \in \dot{M}} F(\underline{v})$$

and

$$(8) \quad \underline{v} \in M, \quad F(\underline{v}) < \mu_2 \Leftrightarrow \underline{v} \in \pm A.$$

By (5), (8),

$$\underline{u} \in \pm A.$$

Take $\underline{u}' = \pm \underline{u}$ so that

$$(9) \quad \underline{u}' \neq \underline{0}, \quad \underline{u}' \in A, \quad F(\underline{u}') < \mu_2 - \mu_1.$$

Let F_A be the linear function which F coincides with on A .

Then

$$\begin{aligned} F_A(\underline{u}') &= F(\underline{u}') \geq F(\underline{a}) && \text{by (7), since } \underline{u}' \neq \underline{0}, \\ &= F_A(\underline{a}). \end{aligned}$$

Hence

$$(10) \quad 0 \leq F_A(\underline{u}' - \underline{a}) < F_A(\underline{u}').$$

Now (5), (7), (8), (10) imply

$$(11) \quad \underline{u}' - \underline{a} \in A, \quad F(\underline{u}' - \underline{a}) < F(\underline{u}') < \mu_2 - \mu_1.$$

Thus (9) implies (11). Repeated use of this fact shows that \underline{u}' is a positive integral multiple of \underline{a} . This establishes (6), and completes the proof of the lemma.

Lemma 7D.2. Let M be a 3-dimensional lattice in R^3 . Let $\underline{u}, \underline{v}$ be linearly independent points in M , and let \underline{f} be a point of R^3 such that

$$(12) \quad \begin{cases} F(\underline{u}) = \alpha, & F(\underline{v}) = \beta, & F(\underline{f}) = \epsilon, \\ \underline{f} \notin \text{lin}(\underline{u}, \underline{v}). \end{cases}$$

Let k be a number such that

$$(13) \quad F(\underline{x}) \geq k \quad \text{for all } \underline{x} \text{ in } M \sim \text{lin}(\underline{u}, \underline{v}).$$

Let P be the plane

$$(14) \quad P = \underline{f} + \text{lin}(\underline{u}, \underline{v}).$$

Suppose that

$$M \cap P \neq \emptyset.$$

Then

$$(15) \quad \epsilon \geq \frac{1}{4} (3k - (\alpha + \beta)).$$

Proof. Let \underline{w} be a point in $M \cap P$. Replacing \underline{w} by $\underline{w} - \underline{u}$, $\underline{w} - \underline{v}$, or $\underline{w} - \underline{u} - \underline{v}$ if necessary, and then replacing \underline{u} by $-\underline{u}$ and/or \underline{v} by $-\underline{v}$ if necessary, we may suppose that

$$\underline{w} = \underline{f} + \lambda \underline{u} + \mu \underline{v},$$

where

$$0 \leq \lambda \leq \frac{1}{2}, \quad 0 \leq \mu \leq \frac{1}{2}.$$

Applying (13) to the point \underline{w} , and to the point

$$2\underline{w} - \underline{u} - \underline{v} = 2\underline{f} + (2\lambda - 1)\underline{u} + (2\mu - 1)\underline{v},$$

and using the convexity of F , we obtain

$$\begin{cases} \epsilon + \lambda\alpha + \mu\beta \geq k \\ 2\epsilon + (1-2\lambda)\alpha + (1-2\mu)\beta \geq k, \end{cases}$$

which immediately gives (15).

Theorem 7D.3. Let M be a 3-dimensional lattice in R^3 , and let \underline{a} , \underline{b} be minimal points as in Lemma 7D.1. Suppose (changing scale if necessary) that

$$(16) \quad F(\underline{a}) = 1 (= \mu_1), \quad F(\underline{b}) = m (= \mu_2).$$

Let P be a lattice plane parallel to $\text{lin}(\underline{a}, \underline{b})$, but not through $\underline{0}$. Then with notation as in (1), (2) (re-naming coordinates if necessary to satisfy (1)), the distance from P to $\underline{0}$ with respect to F satisfies

$$(17) \quad \tau \geq \max \left\{ \frac{m-2}{2}, \frac{2-m}{4}, \frac{1}{6m} \right\}$$

Proof. Apply Lemma 7D.2 with

$$\underline{u} = \underline{a}, \quad \underline{v} = \underline{b}, \quad \underline{f} = (t/\gamma_1, 0, 0).$$

Since (2) states that

$$\tau = |t/\gamma_1|,$$

the lemma gives

$$(18) \quad \tau \geq \frac{1}{4}(3k - (1+m)).$$

By Lemma 7D.1, condition (13) of Lemma 7D.2 will be satisfied if $k = m-1$ or $k = 1$. Substituting these values of k into (18) gives the first two inequalities of (17).

The third inequality in (17) is obtained by a simple application of Minkowski's convex body theorem. The parallel-piped determined by \underline{a} , \underline{b} , \underline{f} contains a fundamental domain for M . Hence

$$\begin{aligned} 1.m.\tau &= F(\underline{a})F(\underline{b})F(\underline{f}) \\ &\geq |\underline{a}||\underline{b}||\underline{f}| \\ &\geq \det M \\ &\geq 2^{-3}. \text{ (Volume of } F(\underline{x}) < 1 \text{), by Minkowski} \\ &= \frac{1}{6}, \text{ as required.} \end{aligned}$$

Note 1. The third inequality in (17), and its proof, are valid if $\underline{a}, \underline{b}$ are arbitrary linearly independent points of M with

$$F(\underline{a}) = 1, \quad F(\underline{b}) = m.$$

Note 2. Application of Lemma 7D.2 to the classical successive minimal points (instead of the minimal points of §7B) and corresponding distance τ' say, yields

$$(19) \quad \tau' \leq \frac{1}{2}\lambda_2 - \frac{1}{4}$$

(from (18) with τ, k, m replaced by $\tau', \lambda_2, \lambda_2$, respectively). Thus, if the minimal point \underline{b} happened to be a second successive minimal point, then (19) would be an improvement on (17). In the case when $k = m = 1$, (17) gives

$$\tau \geq \frac{1}{4}$$

which is easily seen to be best possible, and similarly, if $\lambda_2 = 1 (= \lambda_1)$, then (19) gives $\tau' \geq \frac{1}{4}$ which is best possible.

§7E. Covering domains.

Let L be a 2-dimensional lattice in the plane P through the origin in R^3 . We shall call a set \mathcal{D} a covering domain for L if

$$(1) \quad P = \bigcup_{u \in L} (\mathcal{D} + u),$$

that is, if L is a covering lattice of \mathcal{D} .

In this section, we will find covering domains \mathcal{D} in which F is small, where F is as in §7A(1), and hence derive an auxiliary result needed for the index problem.

Suppose that P has equation

$$(2) \quad \gamma_1 x_1 + \gamma_2 x_2 + \gamma_3 x_3 = 0 ,$$

where

$$(3) \quad |\gamma_1| = \max_i |\gamma_i| > 0 .$$

It follows that if \underline{x} is in P then

$$(4) \quad |x_1| = |(\gamma_2 x_2 + \gamma_3 x_3)/\gamma_1| \leq |x_2| + |x_3| .$$

As explained in §7C, P has four or six coordinate sectors.

We will be concerned with the quadrants $\pm A$, $\pm B$ in P , where

$$(5) \quad \begin{cases} A = \{\underline{x} | \underline{x} \in P, x_2 \geq 0, x_3 \geq 0\} , \\ B = \{\underline{x} | \underline{x} \in P, x_2 \leq 0, x_3 \geq 0\} , \end{cases}$$

which are non-degenerate since $\gamma_1 \neq 0$.

We observe the following immediate consequence of (2), (3) and (4).

Lemma 7E.1. If \underline{x} and \underline{y} are in P , then

$$x_2 = y_2 \quad \text{and} \quad |x_3| < |y_3| \Rightarrow F(\underline{x}) < F(\underline{y}) ,$$

$$x_3 = y_3 \quad \text{and} \quad |x_2| < |y_2| \Rightarrow F(\underline{x}) < F(\underline{y}) .$$

We shall use this observation in the proof of the following lemma.

Lemma 7E.2. Let \underline{u} , \underline{v} be linearly independent points of L such that

$$\underline{u}, \underline{v} \in A, \quad \underline{v} - \underline{u} \in B .$$

Let

$$C = \{\underline{x} | \underline{x} \in P, -u_2 \leq x_2 \leq u_2, 0 \leq x_3 \leq v_3\} ,$$

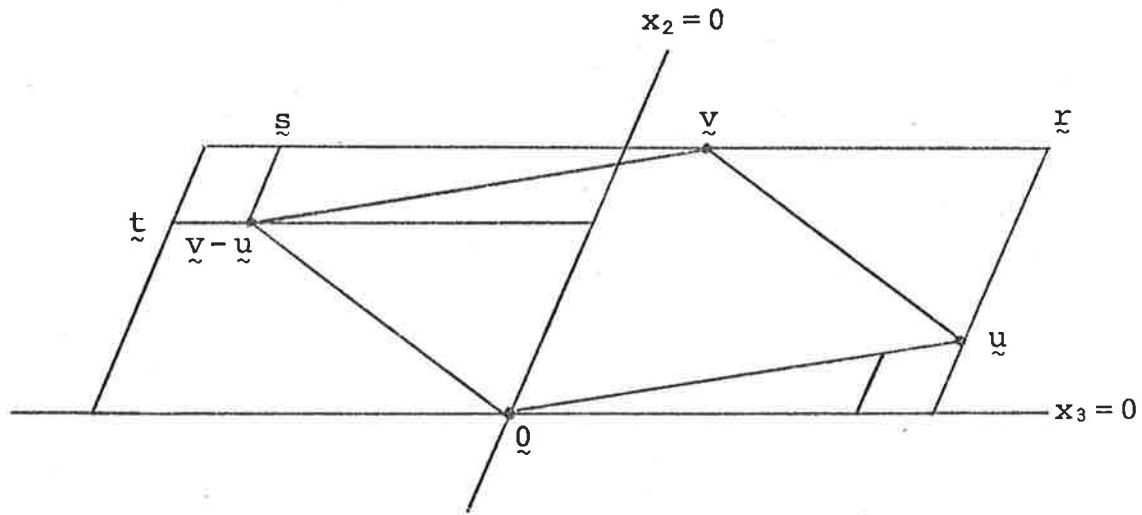
$$\mathcal{D}_1 = C \cap A ,$$

$$\mathcal{D}_2 = C \cap B \sim \{\underline{x} | x_2 < v_2 - u_2, x_3 > v_3 - u_3\} .$$

Then

- (i) \mathcal{D}_1 and \mathcal{D}_2 are covering domains for L , and
- (ii) for all \underline{x} in \mathcal{D}_1 or \mathcal{D}_2

$$F(\underline{x}) \leq \min\{F(\underline{u}) + F(\underline{v}), F(\underline{u}) + F(\underline{v}-\underline{u}), F(\underline{v}) + F(\underline{v}-\underline{u})\}.$$



Proof. (i) The parallelogram F with vertices $0, u, v, v-u$ is certainly a covering domain for L , and if \underline{x} is any point in F then either \underline{x} itself or $\underline{x} + \underline{u}$ or $\underline{x} + \underline{u} - \underline{v}$ is in \mathcal{D}_1 , as illustrated in the diagram. Hence \mathcal{D}_1 is a covering domain for L , and the proof for \mathcal{D}_2 is similar.

(ii) Let $\underline{r}, \underline{s}, \underline{t}$ be the points in P such that

$$r_2 = u_2, r_3 = v_3; \quad s_2 = v_2 - u_2, s_3 = v_3; \quad t_2 = -u_2, t_3 = v_3 - u_3.$$

Then by Lemma 7E.1 the maximum values of F on \mathcal{D}_1 and \mathcal{D}_2 are attained at \underline{r} and \underline{s} or \underline{t} . The result now follows by using Lemma 7E.1 together with the convexity of F . For example,

$$F(\underline{r}) \leq F(\underline{u} + \underline{v}) \leq F(\underline{u}) + F(\underline{v}),$$

$$F(\underline{r}) \leq F(\underline{v}) + F(\underline{r} - \underline{v}) \leq F(\underline{v}) + F(\underline{u} - \underline{v}).$$

We have the following immediate corollary.

Corollary 7E.3. Let L, P, A, B be as in Lemma 7E.2, and let $u, v, v-u$ be points of L not all lying in $\pm A$, and not all lying in $\pm B$. Let δ be any point of P . Then in each quadrant $\pm A, \pm B$ there is a point \tilde{x} such that

$$(6) \quad \tilde{x} \in L + \delta, \quad F(\tilde{x}) \leq F(u) + F(v).$$

We now introduce further notation in the case when there are six sectors in P as in (2), (3), that is, when $\gamma_2\gamma_3 \neq 0$. We then suppose, by mapping (x_1, x_2, x_3) to $(x_1, x_2, -x_3)$ if necessary, that

$$(7) \quad \gamma_2\gamma_3 > 0.$$

In this case the line $x_1 = 0$ in P divides the quadrant B into two sectors, B_1, B_2 , say, so that

$$(8) \quad B = B_1 \cup B_2.$$

For each sector E , let F_E denote the linear function that coincides with F on E . We then have the following lemma, which is easily checked using (4).

Lemma 7E.4. Suppose the plane P given by (2) has six coordinate sectors and (3), (5), (7) and (8) hold. Then

$$F_{B_1}(\tilde{x}) > 0 \quad \text{for all } \tilde{x} \text{ in } B_2,$$

except possibly when \tilde{x} lies on the common boundary of $\pm A$ and $\pm B$, and a similar result holds with B_1 and B_2 interchanged.

The lemma is equivalent to the statement that the lines in P on which F_{B_1}, F_{B_2} vanish must lie in $\pm A$, possibly on its boundary. We now use this lemma together with Corollary 7E.3 to prove the following auxiliary result for our index problem.

Theorem 7E.5. Let L be a 2-dimensional lattice in the plane P in R^3 such that (2), (3) hold. Let \tilde{a}, \tilde{b} , be

minimal points of L as defined in §7C (3), (4) and let

$$F(\underline{a}) = 1, \quad F(\underline{b}) = m.$$

Let A, B be the quadrants in P as defined in (5), and let δ be any point of P . Then in each of the four quadrants $\pm A, \pm B$ there is a point \underline{x} such that

$$\underline{x} \in L + \delta, \quad F(\underline{x}) \leq m + 1.$$

Proof. By Corollary 7E.3, it is sufficient to show that the minimal points $\underline{a}, \underline{b}$ chosen appropriately if not unique, satisfy the condition that $\underline{a}, \underline{b}, \underline{b} - \underline{a}$ do not all lie in $\pm A$, or all in $\pm B$. We will not use the fact that $m \geq 1$, so that $\underline{a}, \underline{b}$ are actually interchangeable in the argument.

Case 1. When $\underline{a} \in \pm A$, and $\underline{b} \in \pm B$ or vice versa, the only way for $\underline{a}, \underline{b}$ to lie in the same pair $\pm A$ or $\pm B$ is for $\underline{a}, \underline{b}$ to lie on the two bounding rays of one quadrant, in which case neither of $\pm(\underline{b} - \underline{a})$ can lie in the same quadrant.

Case 2. When Case 1 does not hold, both \underline{a} and \underline{b} must lie in $\pm B$ but in different sector pairs $\pm B_i$, and we may suppose that $\underline{a} \in B_1, \underline{b} \in B_2$. If $F_{B_2}(\underline{a}) = 0$, we replace \underline{b} by $\underline{b} - h\underline{a}$ where $h \in \mathbb{Z}, h \geq 0$ so that the new \underline{b} still satisfies $\underline{b} \in B_2, F(\underline{b}) = m$, but also $\underline{b} - \underline{a} \in B_2$. If $F_{B_2}(\underline{a}) \neq 0$, then by Lemma 7E.4, $F_{B_2}(\underline{a}) > 0$, so $F_{B_2}(\underline{b} - \underline{a}) < F(\underline{b})$, which would contradict the minimality of $F(\underline{b})$ in $B_2 \sim \ell(\underline{a})$ if we had $\underline{b} - \underline{a} \in B_2$. Hence $\underline{b} - \underline{a} \notin B_2$. Similarly, we may suppose that $\underline{a} - \underline{b} \in B_1$, and so $\underline{b} - \underline{a} \notin -B_1$.

Now if $\underline{b} - \underline{a} \in B_1$, then $\underline{b} = (\underline{b} - \underline{a}) + \underline{a}$ would belong to B_1 , so \underline{b} and hence all three of $\underline{b}, \underline{b} - \underline{a}, \underline{a}$ would lie on the common boundary of B_1 and B_2 , a contradiction since $\underline{b} \notin \ell(\underline{a})$. Thus $\underline{b} - \underline{a} \notin B_1$. Similarly, $\underline{a} - \underline{b} \notin B_2$, so

$\underline{b} - \underline{a} \in -B_2$. Thus, we have that $\underline{b} - \underline{a} \in \pm B_1, \pm B_2$, so $\underline{b} - \underline{a} \in B$, as required.

§7F. Bounds for the index.

Throughout this section, M is a 3-dimensional lattice in R^3 having no non-zero points on the coordinate planes, and F is as in §7A(1). Let

$$(1) \quad X = \{\underline{a}, \underline{b}, \underline{c}, \underline{d}\}, \quad X' = \{\underline{a}, \underline{b}, \underline{c}, \underline{d}'\}$$

be as in §7B (1), (4) with

$$(2) \quad F(\underline{a}) = 1, \quad F(\underline{b}) = m \geq 1.$$

We will use the results of §7D, 7E to obtain bounds for the indices of the sublattices $\ell(X)$ and $\ell(X')$ in M .

Let the plane

$$(3) \quad P = \text{lin}(\underline{a}, \underline{b})$$

have equation

$$(4) \quad \gamma_1 x_1 + \gamma_2 x_2 + \gamma_3 x_3 = 0,$$

where

$$(5) \quad -\gamma_1 \geq \gamma_2 \geq \gamma_3 \geq 0.$$

Condition (5) implies that for all \underline{x} in P

$$(6) \quad |x_1| \leq \frac{\gamma_2}{-\gamma_1 + \gamma_2} F(\underline{x}) \leq \frac{1}{2} F(\underline{x}).$$

Since (5) also implies that (4) of §7E holds, we will be able to use the results obtained there on the quadrants $\pm A$, $\pm B$ of P (bounded by the lines $x_2 = 0$, $x_3 = 0$ in P). Using the notation A_{+++} , etc., for the octants as at the beginning of §7A, we note that

$$(7) \quad P \cap A_{-++} = \{0\},$$

which follows from (5). Since M has no point on the coordinate planes, we have $\gamma_2 \neq 0$. But γ_3 may be zero, in which case we also have $P \cap A_{+-+} = \{0\}$.

For $\xi \in \mathbb{R}$, let $P(\xi)$ denote the plane

$$(8) \quad P(\xi) = (\xi, 0, 0) + P,$$

and recall that by §7D (1), (2),

$$(9) \quad \underline{u} \in P(\xi) \Rightarrow F(\underline{u}) \geq |\xi|.$$

We define

$$(10) \quad \tau_0 = \min\{|\xi| \mid \xi \neq 0, P(\xi) \text{ contains a point of } M\}.$$

Thus, by §7D (1), (2), τ_0 is the distance with respect to F between adjacent lattice planes parallel to P . We recall that by Theorem 7D.3

$$(11) \quad \tau_0 \geq \max\left\{\frac{m-2}{2}, \frac{2-m}{4}, \frac{1}{6m}\right\}.$$

We shall consider the index of $\ell(\underline{a}, \underline{b}, \underline{u})$ in M for certain vectors \underline{u} . By Theorem 7C.2 we have

$$\ell(\underline{a}, \underline{b}) = M \cap P,$$

and it therefore follows that, for any \underline{u} in M and integral h ,

$$(12) \quad \underline{u} \in M \cap P(h\tau_0) \Leftrightarrow (M : \ell(\underline{a}, \underline{b}, \underline{u})) = |h|.$$

The following lemma produces special points in $M \cap P(h\tau_0)$ for suitable h .

Lemma 7F.1. Let $M, \underline{a}, \underline{b}, P, \tau_0$ be as above, satisfying (2) to (5) and (10), and let $h \in \mathbb{Z}$.

(i) If

$$(13) \quad h\tau_0 \geq \frac{m+1}{2}$$

then each of the octants $A_{+++}, A_{+-+}, A_{+--}, A_{++-}$ contains a point \underline{w} such that

$$(14) \quad \begin{cases} \underline{w} \in M \cap P(h\tau_0), \\ F(\underline{w}) \leq h\tau_0 + m + 1. \end{cases}$$

(ii) If $h \geq 0$ the octant A_{+++} contains a point \underline{w} such that

$$\underline{w} = \underline{v} + (h\tau_0, 0, 0) \in M \cap P(h\tau_0),$$

where

$$\underline{v} \in A_{+++}, \quad F(\underline{v}) \leq m + 1.$$

Proof. Let \underline{u} be any point of $M \cap P(h\tau_0)$, and let

$$\underline{\delta} = \underline{u} - (h\tau_0, 0, 0).$$

By applying Theorem 7E.5 with this $\underline{\delta}$, we see that each quadrant $E (= \pm A, \pm B)$ contains a point \underline{v} such that

$$\underline{v} \in E, \quad F(\underline{v}) \leq m + 1,$$

$$\underline{w} = \underline{v} + (h\tau_0, 0, 0) \in M \cap P(h\tau_0).$$

Since, by (6), each \underline{v} has $|v_1| \leq \frac{1}{2}(m+1)$, the condition (13) ensures that each \underline{w} has $w_1 \geq 0$, and this gives the points required in (i). In the case of the octant A_{+++} , (ii) follows from the above and (7).

We can now give a result on indices.

Lemma 7F.2. Let $M, \underline{a}, \underline{b}, P, \tau_0$ be as above, satisfying (2) to (5) and (10). Let G be any octant in R^3 , and let \underline{u} be a point of $G \sim P$ with minimal $F(\underline{x})$. Then

$$(i) \quad (M: \ell(\underline{a}, \underline{b}, \underline{u})) < \frac{3}{2} \frac{m+1}{\tau_0} + 1,$$

and

$$(ii) \quad \text{if } G = \pm A_{+++}, \text{ then}$$

$$(M: \ell(\underline{a}, \underline{b}, \underline{u})) \leq \frac{m+1}{\tau_0} + 1.$$

Proof. Without loss of generality, we may suppose that G is one of the octants $A_{+++}, A_{++-}, A_{+-+}, A_{-++}$. Apply Lemma 7F.1 with h as the least integer greater than or equal to $(m+1)/2\tau_0$, so that

$$(15) \quad 1 \leq h < \frac{m+1}{2\tau_0} + 1,$$

to obtain a point \underline{w} in the same octant as \underline{u} and satisfying (14), so that $\underline{w} \in M \sim P$, and

$$(16) \quad F(\underline{w}) \leq h\tau_0 + (m+1).$$

Let k be the required index, so that, by (12),

$$\underline{u} \in M \cap P(k\tau_0),$$

hence by (9), we have

$$(17) \quad F(\underline{u}) \geq k\tau_0.$$

Suppose that (i) does not hold. Then

$$k\tau_0 \geq \frac{m+1}{2} + \tau_0 + (m+1) > 0,$$

so by (15), (17),

$$F(\underline{u}) > h\tau_0 + (m+1), \quad \underline{u} \notin P.$$

By (16), this contradicts the minimality of \underline{u} . Thus (i) holds.

The result (ii) is proved similarly, using (ii) of Lemma 7F.1 with $h = 1$.

We now combine our results to give our final index estimates.

Theorem 7F.3. Let M be a 3-dimensional lattice in R^3 having no non-zero points on the coordinate planes. Let X, X' be as in (1) above and in §7B (1), (4), and let

$$m = F(\underline{b})/F(\underline{a}) \geq 1,$$

$$\omega = \max \left\{ \frac{m-2}{2}, \frac{2-m}{4}, \frac{1}{6m} \right\}.$$

(i) Each of the indices $\{M:\ell(X)\}, \{M:\ell(X')\}$ is positive and less than

$$\frac{3}{2} \frac{m+1}{\omega} + 1.$$

(ii) Let P be as in (3), (4), (5). If neither \underline{a} nor \underline{b} lies in the quadrant

$$A = \{\underline{x} | \underline{x} \in P, x_2 \geq 0, x_3 \geq 0\},$$

then

$$\{M:\ell(X')\} \leq \frac{m+1}{\omega} + 1.$$

Corollary 7F.4. For all m , the indices in (i) are at most 62 and if $m \geq 11$ these indices are at most 4. For all

m , the index in (ii) is at most 41 and if $m \geq 8$ this index is at most 3.

Note. The above bounds are rather crude for small m . It should be possible to drastically improve them by obtaining an improved version of (17) of Theorem 7D.3, in particular for the case when m is near 2.

Proof of Theorem 7F.3. The indices $(M:l(X))$ and $(M:l(X'))$ are just the greatest common divisors of $(M:l(\underline{a}, \underline{b}, \underline{c}))$ with $(M:l(\underline{a}, \underline{b}, \underline{d}))$ and $(M:l(\underline{a}, \underline{b}, \underline{d}'))$ respectively. Hence to prove (i) it is sufficient to prove that

$$(18) \quad 0 < (M:l(\underline{a}, \underline{b}, \underline{u})) < \frac{3}{2} \frac{m+1}{\omega} + 1$$

holds when $\underline{u} = \underline{c}$, and to prove (ii) it is sufficient to prove that under the assumptions on $\underline{a}, \underline{b}$ in (ii), we have

$$(19) \quad (M:l(\underline{a}, \underline{b}, \underline{u})) < \frac{m+1}{\omega} + 1$$

when \underline{u} is one of $\underline{c}, \underline{d}'$. Now if \underline{u} is any of the points $\pm \underline{c}, \pm \underline{d}'$ and G is the octant containing \underline{u} , then \underline{u} is the minimal point of $G \sim \text{lin}(\underline{a}, \underline{b})$. Assume without loss of generality that P is as in (ii). Then on applying Lemma 7F.2(i) and using (11), we obtain (18), in particular with $\underline{u} = \underline{c}$.

If neither of $\underline{a}, \underline{b}$ lies in A , then $\pm \underline{a}, \pm \underline{b}$ must lie in A_{++-}, A_{+-+} (or vice versa) and one of $\pm \underline{c}, \pm \underline{d}'$ must lie in A_{+++} . Taking this point as \underline{u} in Lemma 7F.2 (ii), and using (11), we obtain (19). This concludes the proof of the theorem.

BIBLIOGRAPHY

- [1] Bantegnie, R., "A propos d'un problème de Mordell sur les octaèdres latticiels", J. London Math. Soc., 37 (1962), 320-328.
- [2] Brown, H. and Mahler, K., "A generalization of Farey sequences: Some exploration via the computer", J. Number Theory 3 (1971), 364-370.
- [3] Cassels, J.W.S., "An introduction to the geometry of numbers", Springer-Verlag, Berlin-Göttingen-Heidelberg, 1959.
- [4] Gaskell, R.W., Klamkin, M.S. and Watson, P., "Triangulations and Pick's Theorem", Math. Mag. 49 (1976), 35-37.
- [5] Grünbaum, B., "Convex polytopes", Interscience, London, 1967.
- [6] Hardy, G.H. and Wright, E.M., "An introduction to the theory of numbers", Clarendon Press, Oxford, 1960.
- [7] Lang, S., "Introduction to diophantine approximations", Addison-Wesley, Reading, Mass., 1966.
- [8] Lekkerkerker, Cornelis Gerrit, "Geometry of numbers", Wolters-Noordhoff, Groningen; North-Holland, Amsterdam-London, 1969.
- [9] Low, L., "A problem of Schinzel on lattice points", Acta Arith. 31 (1976), 385-388.
- [10] Mordell, L.J., "Lattice octahedra", Canad. J. Math. 12 (1960), 297-302.
- [11] Schinzel, A., "On the reducibility of polynomials and in particular of trinomials", Acta Arith. 11 (1965), 1-34.
- [12] Schmidt, W.M., "A problem of Schinzel on lattice points", Acta Arith. 15 (1969), 199-203.