11:44 a.m. April 2, 2011

## Hermite and Smith forms

Bill Casselman
University of British Columbia
cass@math.ubc.ca

Integral matrices may be reduced to simple form through left and right multiplication by invertible integral matrices. This has many consequences, among them the fact that every finite abelian group is a product of cyclic groups, and it is also a key fact in Siegel's computation of the volume of $\mathrm{SL}_n(\mathbb{Z})\backslash X_n$, where $X_n$ is the space of all positive definite matrices of determinant 1.

This is a special if important case of a general result about principal ideal domains, but worthwhile covering separately because proofs are accompanied by explicit algorithms.

**Contents**

### 1. Lattice arithmetic

There are several lattice computations necessary to prove the Theorem constructively. The first result we'll need is a basic tool that will be used several times:

[euclid-alg] **Lemma 1.1.** (Euclidean algorithm) *Given an integral vector $v = [v_1 \ \ldots \ v_n]$, one can find an invertible integral matrix $A$ such that*

$$vA = [0\ 0\ \ldots\ d]$$

*where $d$ is (necessarily) the greatest common divisor of the coordinates of $v$.*

*Proof.* The case $n = 1$ is trivial. Let's look next at the case $n = 2$. Given the vector $[m, n]$, we carry out the Euclidean algorithm to find the greatest common divisor $d$ of $m$ and $n$, keeping track of a few extra things as we go.

Precisely, we start with the vector $[m_0, n_0] = [m, n]$ and the matrix $A_0 = I$, the $2 \times 2$ identity. We are going to calculate a sequence of vectors and matrices $[m_i, n_i]$ and $A_i$, satisfying at all times the equation $[m_0, n_0]A_i = [m_i, n_i]$, and winding up with some $[m_i, n_i] = [0, d]$. As long as $m_i \neq 0$, in going from $i$ to $i + 1$ we divide $n_i$ by $m_i$:

$$n_i = qm_i + r$$
$$n_{i+1} = m_i$$
$$m_{i+1} = r$$
$$A_{i+1} = A_i \begin{bmatrix} -q & 1 \\ 1 & 0 \end{bmatrix}$$

This works since

$$[m_{i+1}\ n_{i+1}] = [(n_i - qm_i)\ m_i] = [m_i\ n_i] \begin{bmatrix} -q & 1 \\ 1 & 0 \end{bmatrix} = [m_0\ n_0]A_i \begin{bmatrix} -q & 1 \\ 1 & 0 \end{bmatrix} .$$

The process stops when division is exact and hence $r = 0$.

Now we continue by induction. We start with

$$v = [v_1 \ v_2 \ \ldots \ v_n] .$$

We may assume we have found an $(n-1) \times (n-1)$ $A_{n-1}$ matrix embedded into an $n \times n$ matrix such that

$$v \begin{bmatrix} A_{n-1} & 0 \\ 0 & 1 \end{bmatrix} = [0 \ 0 \ \ldots d_{n-1} \ v_n]$$

where $d_{n-1}$ is the gcd of the first $n-1$ coordinates of $v$. But then according to the case $n=2$ we can multiply by an embedded $2 \times 2$ matrix to get $[0 \ 0 \ \ldots \ d_n]$. ∎

**[subsat] Lemma 1.2.** (Subspace saturation) *If $U$ is a rational vector subspace of $\mathbb{Q}^n$ then the intersection $U \cap \mathbb{Z}^n$ is a summand of $\mathbb{Z}^n$.*

I'll assume that $U$ is given in the form of a basis of $m$ linearly independent rational vectors, put as columns into a matrix I am afraid I shall call $U$ from now on. The proof will explain how to find a $\mathbb{Z}$-basis of the intersection as well as a complement in $\mathbb{Z}^n$.

For $m \leq n$ define $I_{n,m}$ to be the $n \times m$ matrix

$$I_{n,m} = \begin{bmatrix} I_m \\ 0 \end{bmatrix} .$$

For $i \leq n$ let $e_{i,n}$ be the $n$-dimensional vector with the $i$-th coordinate 1 and all others 0.

**[exsat] Lemma 1.3.** (Explicit saturation) *Suppose $U$ to be an $n \times m$ matrix whose columns are linearly independent rational vectors in $\mathbb{Q}^n$. We can find a matrix $A$ in $\mathrm{GL}_n(\mathbb{Z})$ and a matrix $B$ in $\mathrm{GL}_m(\mathbb{Q})$ such that $AUB$ is $I_{n,m}$.*

Under the hypothesis of independence, of course $m \leq n$. I recall that $\mathrm{GL}_n(\mathbb{Z})$ is the group of integral matrices of determinant $\pm 1$, which is to say those integral matrices with integral inverses.

Why does this imply the Lemma? The matrix $UB$ is a new rational basis of the vector space spanned by the columns of $U$. Since $UB = A^{-1}I_{n,m}$, it is also the first $m$ columns of the matrix $A^{-1}$, whose columns make up a basis of $\mathbb{Z}^n$. This is exactly what we want.

*Proof* I shall tell exactly how to get $A$ and $B$. The proof proceeds by induction on $m$, and the case $m=1$ is a simple variant of the previous Lemma.

The case $m=1$ of this asserts that if $u$ is any vector in $\mathbb{Q}^n$, there exists some $b$ in $\mathbb{Q}^{\times}$ and an $A$ in $\mathrm{GL}_n(\mathbb{Z})$ such that $Aub$ is the column vector $v$ with $v_1 = 1$, $v_i = 0$ for $i > 1$. We start by multiplying $u$ by some integer $p$ to make $pu$ itself integral. The Lemma in slight disguise now finds $A$ such that $v = Apu$ has all $v_i = 0$ for all $i > 1$, and $v_1 = q$ where $q$ is the gcd of the coefficients $pu_i$. That is to say, $v = qe_{1,n}$. But then $(p/q)u$ is still integral, and $A(p/q)u = e_{1,n}$.

Now suppose $m > 1$, and assume the Lemma to be true for $m-1$. We can find $A_*$ and $B_*$ such that $A_* u B_*$ is a matrix whose first $m-1$ columns are $I_{m-1,n}$:

$$A_* U B_* = \begin{bmatrix} I_{m-1} & u_{m-1} \\ 0 & u_{n-m+1} \end{bmatrix} .$$

Here $u_{m-1}$ is a column vector of length $m-1$, and $u_{n-m+1}$ one of length $n-m+1$. Elementary column operations, amounting to multiplication of this on the right by certain triangular matrices, will make $u_{m-1} = 0$, and then we can apply the case $m=1$ to get $u_{n-m+1} = e_{m,n}$. ∎

The proof tells you how to calculate $A$ and $B$ as you go along, but it also tells you how to calculate $A^{-1}$ and $B^{-1}$, since, for example, multiplying $A$ on the left by an embedded $2 \times 2$ matrix $S$ is no easier than multiplying $A^{-1}$ on the right by $S^{-1}$, which is trivial to compute.

## 2. Hermite

I'll say that an integral matrix $E$ is in **Hermite normal form** if it looks like

$$\begin{bmatrix} 0 & H \end{bmatrix},$$

where $H$ is a matrix satisfying certain echelon conditions. Suppose it has $d$ columns. We require first of all that (a) no column of $H$ is all zeroes. Then in each of the columns $j$ there exists a last non-zero entry, say in row $r(j)$. (b) The entry $p(j) = h_{r(j),j}$ is positive. These are called the pivot entries. (c) If $j < k$ then $r(j) < r(k)$. (d) If $k > j$ then $h_{r(j),k}$ is in the range $[0, p(j))$. The general shape of a matrix in Hermite normal form is thus something like

$$\begin{bmatrix} \star & \star & \star \\ \bullet & \diamond & \diamond \\ 0 & \star & \star \\ 0 & \bullet & \diamond \\ 0 & 0 & \bullet \\ 0 & 0 & 0 \end{bmatrix},$$

where $\bullet \neq 0$, $\star$ is arbitrary, and $\diamond$ is subject to the range condition (d).

**[hermite] Theorem 2.1.** (Hermite normal form) *If $M$ is any integral $r \times c$ matrix, one can find $B$ in $\mathrm{GL}_c(\mathbb{Z})$ such that $MB = H$ is in Hermite normal form. This normal form is unique.*

The non-zero columns of $H$ make up a distinguished basis of the lattice spanned by the columns of $M$, relative to the standard flag

$$0 \subset \mathbb{Q} \subset \mathbb{Q}^2 \subset \ldots \subset \mathbb{Q}^n.$$

*Proof* Finding the Hermite normal form of a matrix requires first of all several applications, row by row from the bottom up, of the Euclidean algorithm. This gives conditions (a), (b), and (c). Condition (d) can then be satisfied by applying some elementary column operations.

Now for uniqueness. The claim is that an $n \times m$ matrix $A$ in Hermite normal form whose columns form a basis of a $\mathbb{Z}$-subgroup $L$ of $\mathbb{Z}^n$ is determined by $L$. The number of its non-zero columns is the rank of $L$, hence certainly determined by $L$.

For each $i$, let $M_i$ be the intersection of $L$ with the subspace $x_k = 0$ for $k > i$. Let $r_j$ be the number of trailing zeroes in column $j$ of $A$. Then $r_1$ is the largest $r$ such that $L \cap M_r \neq 0$, and more generally $r_j$ si the largest $r$ such that $L \cap M_r$ has rank $j$. Therefore the $r_j$ are determined by $L$ alone. The submodule $L_j$ spanned by the first $j$ columns of $A$ is the intersection of $L$ with $L_{r_j}$, hence also determined by $L$ alone.

If $L$ has rank one, then the column $\ell_1$ of $A$ is a basis of $L$. It is the unique basis of $L$ with last entry positive, which proves the claim for rank one.

Fome here we go by induction on the rank $c$ of $L$. The submodule $L_{c-1}$ spanned by the first $c-1$ columns of $A$ is uniquely determined by $L$, and we may apply induction to it. Therefore the first $c-1$ columns of $A$ are uniquely determined by $A$. The last column of $A$ is a basis elemnt of $L/L_{c-1}$, to which we may apply induction. Therefore it is uniquely determined modulo $L_{c-1}$. But then condition (d) determines it uniquely. ∎

**[eisenstein] Corollary 2.2.** *Every non-singular matrix in $M_n(\mathbb{Z})$ with positive determinant is equivalent modulo right multiplication by a matrix in $\mathrm{SL}_n(\mathbb{Z})$ to a unique matrix in Hermitian normal form with positive diagonal entries.*

As a consequence:

**[siegel-count] Proposition 2.3.** *The number of integral $n \times n$ matrices with determinant $d > 0$ modulo right multiplication by elements of $\mathrm{SL}_n(\mathbb{Z})$ is equal to*

$$\sum d_1^{n-1} d_2^{n-2} \ldots d_1^0$$

*with the sum over all $(d_i)$ with $\prod d_i = d$.*

In Siegel's computation of the volume of $\mathrm{SL}_n(\mathbb{Z}) \backslash \mathrm{SL}_n(\mathbb{Z})$ this is attributed to Eisenstein, but I have not been able to locate a precise reference among his collected works.

### 3. Smith

A matrix is said to be in **Smith normal form** (presumably named after the prominent nineteenth century English mathematician H. J. Stephen Smith) if it looks likes

$$\begin{bmatrix} 0 & D \\ 0 & 0 \end{bmatrix},$$

where $D$ is a diagonal matrix with positive entries such that $d_{i,i} | d_{i+1,i+1}$.

**[smith] Lemma 3.1.** (Smith normal form) *If $M$ is an integral matrix of size $r \times c$, we can find $A$ in $\mathrm{GL}_r(\mathbb{Z})$ and $B$ in $\mathrm{GL}_c(\mathbb{Z})$ such that $A\,M\,B = S$ is in Smith normal form.*

Here, $MB$ is a basis of the lattice $L_M$ generated by the columns of $M$. If $S$ has $k$ non-zero columns, then since $MB = A^{-1}S$, the first $k$ columns $a_i$ of $A^{-1}$ are part of a basis of $\mathbb{Z}^n$ such that $d_{i,i}a_i$ make up a basis of $L_M$. The matrix $B^{-1}$ expresses the columns of $M$ in terms of that basis. As earlier, it will be no more difficult to find $A^{-1}$ and $B^{-1}$ than $A$ and $B$.

Thus finding the Smith normal form is the same as implementing the principal divisor theorem.

*Proof* First put the matrix in Hermite normal form. Then multiply it on the left by a matrix in $\mathrm{GL}_r(\mathbb{Z})$ to get it in the form

$$\begin{bmatrix} 0 & d & * & \ldots & * \\ 0 & 0 & * & \ldots & * \\ \ldots & & & & \\ 0 & 0 & * & \ldots & * \end{bmatrix}.$$

There are now two possibilities: (i) The corner entry $d$ is the gcd of the first row. We can tell whether this is true by running along the first row, applying an elementary column operation to replace an entry by its remainder upon division by $d$. If these remainders were all 0, our matrix looks like

$$\begin{bmatrix} 0 & d_{1,1} & 0 & \ldots & 0 \\ 0 & 0 & * & \ldots & * \\ \ldots & & & & \\ 0 & 0 & * & \ldots & * \end{bmatrix}.$$

We move on to the next column to get $d_{2,2}$. Etc. (ii) The entry $d$ is not the gcd of the top row, and some of those remainders were not zero. In this case, we apply the Euclidean algorithm to replace $d$ by the gcd of the row. This may, however, place some non-zero integers in the first column, so we have to go back to the start. We keep on applying the Euclidean algorithm to the first column and row, but in each cycle the corner entry decreases, so we must eventually break the loop.

At the end of this part of the computation, we'll have a diagonal matrix $D$, but one which might not satisfy the divisibility condition. To obtain that, we perform several times an operation essentially in

$\mathrm{GL}_2(\mathbb{Z})$. This operation, in effect, deals with a special case of our problem. Suppose we are given a diagonal integral matrix

$$\begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix}.$$

We want to multiply on left and right by matrices in $\mathrm{GL}_2(\mathbb{Z})$ to get a similar diagonal matrix, but with $a|b$. Let $d$ be the gcd of $a$ and $b$. Perform the usual Euclidean algorithm to find $k$ and $\ell$ such that $ka + \ell b = d$. Adding $k$ times the first column to the second gives

$$\begin{bmatrix} a & ka \\ 0 & b \end{bmatrix},$$

and then adding $\ell$ times the second row to the first gives

$$\begin{bmatrix} a & ka + \ell b \\ 0 & b \end{bmatrix} = \begin{bmatrix} a & d \\ 0 & b \end{bmatrix}.$$

A signed swap of columns gives

$$\begin{bmatrix} d & -a \\ b & 0 \end{bmatrix}.$$

Since $d$ divides both $a$ and $b$ we can subtract $b/d$ times the first row from the second to get

$$\begin{bmatrix} d & -a \\ 0 & ab/d \end{bmatrix};$$

and finally add $a/d$ times the first column to the second to get

$$\begin{bmatrix} d & 0 \\ 0 & ab/d \end{bmatrix}.$$

We apply this operation to a large diagonal matrix several times, to all pairs $d_{i,i}$ and $d_{j,j}$ with $i < j$. First we get $d_{1,1}$ to divide all $d_{j,j}$ with $j > 1$, then get $d_{2,2}$ to divide all $d_{j,j}$ with $j > 2$, etc. ∎

Chapter 3 of [Pohst-Zassenhaus:1989] explains in more detail how to compute Hermite and Smith normal forms of an arbitrary integral matrix. In particular, they point out the difficulties that arise, which do not appear in my account. Neither of the two normal forms is usually an easy computation, but finding the Hermite form is noticeably simpler, and it is in any event a good first step towards the Smith form. Among other things it tells us immediately the column rank at hand. In one situation we'll be dealing with, our matrix will start off in a particularly good Hermite form.

## 4. References

**1.** M. Pohst and H. Zassenhaus, **Algorithmic algebraic number theory**, Cambridge University Press, 1989.

**2.** Carl Ludwig Siegel, (1936),