

Theorem 38. Let R be a ring with 1 and let M be an R -module. Then M is contained in an injective R -module.

Proof: A proof is outlined in Exercises 15 to 17.

It is possible to prove a sharper result than Theorem 38, namely that there is a *minimal* injective R -module H containing M in the sense that any injective map of M into an injective R -module Q factors through H . More precisely, if $M \subseteq Q$ for an injective R -module Q then there is an injection $\iota : H \hookrightarrow Q$ that restricts to the identity map on M ; using ι to identify H as a subset of Q we have $M \subseteq H \subseteq Q$. (cf. Theorem 57.13 in *Representation Theory of Finite Groups and Associative Algebras* by C. Curtis and I. Reiner, John Wiley & Sons, 1966). This module H is called the *injective hull* or *injective envelope* of M . The universal property of the injective hull of M with respect to inclusions of M into injective R -modules should be compared to the universal property with respect to homomorphisms of M of the free module $F(A)$ on a set of generators A for M in Theorem 6. For example, the injective hull of \mathbb{Z} is \mathbb{Q} , and the injective hull of any field is itself (cf. the exercises).

Flat Modules and $D \otimes_R \underline{}$

We now consider the behavior of extensions $0 \longrightarrow L \xrightarrow{\psi} M \xrightarrow{\varphi} N \longrightarrow 0$ of R -modules with respect to tensor products.

Suppose that D is a *right* R -module. For any homomorphism $f : X \rightarrow Y$ of left R -modules we obtain a homomorphism $1 \otimes f : D \otimes_R X \rightarrow D \otimes_R Y$ of abelian groups (Theorem 13). If in addition D is an (S, R) -bimodule (for example, when $S = R$ is commutative and D is given the standard (R, R) -bimodule structure as in Section 4), then $1 \otimes f$ is a homomorphism of left S -modules. Put another way,

$$D \otimes_R \underline{} : X \longrightarrow D \otimes_R X$$

is a *covariant functor* from the category of left R -modules to the category of abelian groups (respectively, to the category of left S -modules when D is an (S, R) -bimodule), cf. Appendix II. In a similar way, if D is a left R -module then $\underline{} \otimes_R D$ is a covariant functor from the category of right R -modules to the category of abelian groups (respectively, to the category of right S -modules when D is an (R, S) -bimodule). Note that, unlike Hom, the tensor product is covariant in both variables, and we shall therefore concentrate on $D \otimes_R \underline{}$, leaving as an exercise the minor alterations necessary for $\underline{} \otimes_R D$.

We have already seen examples where the map $1 \otimes \psi : D \otimes_R L \rightarrow D \otimes_R M$ induced by an injective map $\psi : L \hookrightarrow M$ is no longer injective (for example the injection $\mathbb{Z} \hookrightarrow \mathbb{Q}$ of \mathbb{Z} -modules induces the zero map from $\mathbb{Z}/2\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Z} = \mathbb{Z}/2\mathbb{Z}$ to $\mathbb{Z}/2\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Q} = 0$). On the other hand, suppose that $\varphi : M \rightarrow N$ is a surjective R -module homomorphism. The tensor product $D \otimes_R N$ is generated as an abelian group by the simple tensors $d \otimes n$ for $d \in D$ and $n \in N$. The surjectivity of φ implies that $n = \varphi(m)$ for some $m \in M$, and then $1 \otimes \varphi(d \otimes m) = d \otimes \varphi(m) = d \otimes n$ shows that $1 \otimes \varphi$ is a surjective homomorphism of abelian groups from $D \otimes_R M$ to $D \otimes_R N$. This proves most of the following theorem.

Theorem 39. Suppose that D is a right R -module and that L , M and N are left R -modules. If

$$0 \longrightarrow L \xrightarrow{\psi} M \xrightarrow{\varphi} N \longrightarrow 0 \quad \text{is exact,}$$

then the associated sequence of abelian groups

$$D \otimes_R L \xrightarrow{1 \otimes \psi} D \otimes_R M \xrightarrow{1 \otimes \varphi} D \otimes_R N \longrightarrow 0 \quad \text{is exact.} \quad (10.13)$$

If D is an (S, R) -bimodule then (13) is an exact sequence of left S -modules. In particular, if $S = R$ is a commutative ring, then (13) is an exact sequence of R -modules with respect to the standard R -module structures. The map $1 \otimes \varphi$ is not in general injective, i.e., the sequence (13) cannot in general be extended to a short exact sequence.

The sequence (13) is exact for all right R -modules D if and only if

$$L \xrightarrow{\psi} M \xrightarrow{\varphi} N \longrightarrow 0 \quad \text{is exact.}$$

Proof: For the first statement it remains to prove the exactness of (13) at $D \otimes_R M$. Since $\varphi \circ \psi = 0$, we have

$$(1 \otimes \varphi) \left(\sum d_i \otimes \psi(l_i) \right) = \sum d_i \otimes (\varphi \circ \psi(l_i)) = 0$$

and it follows that $\text{image}(1 \otimes \psi) \subseteq \ker(1 \otimes \varphi)$. In particular, there is a natural projection $\pi : (D \otimes_R M)/\text{image}(1 \otimes \psi) \rightarrow (D \otimes_R M)/\ker(1 \otimes \varphi) = D \otimes_R N$. The composite of the two projection homomorphisms

$$D \otimes_R M \rightarrow (D \otimes_R M)/\text{image}(1 \otimes \psi) \xrightarrow{\pi} D \otimes_R N$$

is the quotient of $D \otimes_R M$ by $\ker(1 \otimes \varphi)$, so is just the map $1 \otimes \varphi$. We shall show that π is an isomorphism, which will show that the kernel of $1 \otimes \varphi$ is just the kernel of the first projection above, i.e., $\text{image}(1 \otimes \psi)$, giving the exactness of (13) at $D \otimes_R M$. To see that π is an isomorphism we define an inverse map. First define $\pi' : D \times N \rightarrow (D \otimes_R M)/\text{image}(1 \otimes \psi)$ by $\pi'((d, n)) = d \otimes m$ for any $m \in M$ with $\varphi(m) = n$. Note that this is well defined: any other element $m' \in M$ mapping to n differs from m by an element in $\ker \varphi = \text{image } \psi$, i.e., $m' = m + \psi(l)$ for some $l \in L$, and $d \otimes \psi(l) \in \text{image}(1 \otimes \psi)$. It is easy to check that π' is a balanced map, so induces a homomorphism $\tilde{\pi} : D \times N \rightarrow (D \otimes_R M)/\text{image}(1 \otimes \psi)$ with $\tilde{\pi}(d \otimes n) = d \otimes m$. Then $\tilde{\pi} \circ \pi(d \otimes m) = \tilde{\pi}(d \otimes \varphi(m)) = d \otimes m$ shows that $\tilde{\pi} \circ \pi = 1$. Similarly, $\pi \circ \tilde{\pi} = 1$, so that π and $\tilde{\pi}$ are inverse isomorphisms, completing the proof that (13) is exact. Note also that the injectivity of ψ was not required for the proof.

Finally, suppose (13) is exact for every right R -module D . In general, $R \otimes_R X \cong X$ for any left R -module X (Example 1 following Corollary 9). Taking $D = R$ the exactness of the sequence $L \xrightarrow{\psi} M \xrightarrow{\varphi} N \longrightarrow 0$ follows.

By Theorem 39, the sequence

$$0 \longrightarrow D \otimes_R L \xrightarrow{1 \otimes \psi} D \otimes_R M \xrightarrow{1 \otimes \varphi} D \otimes_R N \longrightarrow 0$$

is not in general exact since $1 \otimes \psi$ need not be injective. If $0 \rightarrow L \xrightarrow{\psi} M \xrightarrow{\varphi} N \rightarrow 0$ is a *split* short exact sequence, however, then since tensor products commute with direct sums by Theorem 17, it follows that

$$0 \longrightarrow D \otimes_R L \xrightarrow{1 \otimes \psi} D \otimes_R M \xrightarrow{1 \otimes \varphi} D \otimes_R N \longrightarrow 0$$

is also a split short exact sequence.

The following result relating to modules D having the property that (13) can always be extended to a short exact sequence is immediate from Theorem 39:

Proposition 40. Let A be a right R -module. Then the following are equivalent:

- (1) For any left R -modules L , M , and N , if

$$0 \longrightarrow L \xrightarrow{\psi} M \xrightarrow{\varphi} N \longrightarrow 0$$

is a short exact sequence, then

$$0 \longrightarrow A \otimes_R L \xrightarrow{1 \otimes \psi} A \otimes_R M \xrightarrow{1 \otimes \varphi} A \otimes_R N \longrightarrow 0$$

is also a short exact sequence.

- (2) For any left R -modules L and M , if $0 \rightarrow L \xrightarrow{\psi} M$ is an exact sequence of left R -modules (i.e., $\psi : L \rightarrow M$ is injective) then $0 \rightarrow A \otimes_R L \xrightarrow{1 \otimes \psi} A \otimes_R M$ is an exact sequence of abelian groups (i.e., $1 \otimes \psi : A \otimes_R L \rightarrow A \otimes_R M$ is injective).

Definition. A right R -module A is called *flat* if it satisfies either of the two equivalent conditions of Proposition 40.

For a fixed right R -module D , the first part of Theorem 39 is referred to by saying that the functor $D \otimes_R \underline{}$ is *right exact*.

Corollary 41. If D is a right R -module, then the functor $D \otimes_R \underline{}$ from the category of left R -modules to the category of abelian groups is right exact. If D is an (S, R) -bimodule (for example when $S = R$ is commutative and D is given the standard R -module structure), then $D \otimes_R \underline{}$ is a right exact functor from the category of left R -modules to the category of left S -modules. The functor is exact if and only if D is a flat R -module.

We have already seen some flat modules:

Corollary 42. Free modules are flat; more generally, projective modules are flat.

Proof: To show that the free R -module F is flat it suffices to show that for any injective map $\psi : L \rightarrow M$ of R -modules L and M the induced map $1 \otimes \psi : F \otimes_R L \rightarrow F \otimes_R M$ is also injective. Suppose first that $F \cong R^n$ is a finitely generated free R -module. In this case $F \otimes_R L = R^n \otimes_R L \cong L^n$ since $R \otimes_R L \cong L$ and tensor products commute with direct sums. Similarly $F \otimes_R M \cong M^n$ and under these isomorphisms

the map $1 \otimes \psi : F \otimes_R L \rightarrow F \otimes_R M$ is just the natural map of L^n to M^n induced by the inclusion ψ in each component. In particular, $1 \otimes \psi$ is injective and it follows that any finitely generated free module is flat. Suppose now that F is an arbitrary free module and that the element $\sum f_i \otimes l_i \in F \otimes_R L$ is mapped to 0 by $1 \otimes \psi$. This means that the element $\sum (f_i, \psi(l_i))$ can be written as a sum of generators as in equation (6) in the previous section in the free group on $F \times M$. Since this sum of elements is finite, all of the first coordinates of the resulting equation lie in some finitely generated free submodule F' of F . Then this equation implies that $\sum f_i \otimes l_i \in F' \otimes_R L$ is mapped to 0 in $F' \otimes_R M$. Since F' is a finitely generated free module, the injectivity we proved above shows that $\sum f_i \otimes l_i$ is 0 in $F' \otimes_R L$ and so also in $F \otimes_R L$. It follows that $1 \otimes \psi$ is injective and hence that F is flat.

Suppose now that P is a projective module. Then P is a direct summand of a free module F (Proposition 30), say $F = P \oplus P'$. If $\psi : L \rightarrow M$ is injective then $1 \otimes \psi : F \otimes_R L \rightarrow F \otimes_R M$ is also injective by what we have already shown. Since $F = P \oplus P'$ and tensor products commute with direct sums, this shows that

$$1 \otimes \psi : (P \otimes_R L) \oplus (P' \otimes_R L) \rightarrow (P \otimes_R M) \oplus (P' \otimes_R M)$$

is injective. Hence $1 \otimes \psi : P \otimes_R L \rightarrow P \otimes_R M$ is injective, proving that P is flat.

Examples

- (1) Since \mathbb{Z} is a projective \mathbb{Z} -module it is flat. The example before Theorem 39 shows that $\mathbb{Z}/2\mathbb{Z}$ not a flat \mathbb{Z} -module.
- (2) The \mathbb{Z} -module \mathbb{Q} is a flat \mathbb{Z} -module, as follows. Suppose $\psi : L \rightarrow M$ is an injective map of \mathbb{Z} -modules. Every element of $\mathbb{Q} \otimes_{\mathbb{Z}} L$ can be written in the form $(1/d) \otimes l$ for some nonzero integer d and some $l \in L$ (Exercise 7 in Section 4). If $(1/d) \otimes l$ is in the kernel of $1 \otimes \psi$ then $(1/d) \otimes \psi(l)$ is 0 in $\mathbb{Q} \otimes_{\mathbb{Z}} M$. By Exercise 8 in Section 4 this means $c\psi(l) = 0$ in M for some nonzero integer c . Then $\psi(c \cdot l) = 0$, and the injectivity of ψ implies $c \cdot l = 0$ in L . But this implies that $(1/d) \otimes l = (1/cd) \otimes (c \cdot l) = 0$ in L , which shows that $1 \otimes \psi$ is injective.
- (3) The \mathbb{Z} -module \mathbb{Q}/\mathbb{Z} is injective (by Proposition 36), but is not flat: the injective map $\psi(z) = 2z$ from \mathbb{Z} to \mathbb{Z} does not remain injective after tensoring with \mathbb{Q}/\mathbb{Z} . $(1 \otimes \psi : \mathbb{Q}/\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Z} \rightarrow \mathbb{Q}/\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Z}$ has the nonzero element $(\frac{1}{2} + \mathbb{Z}) \otimes 1$ in its kernel — identifying $\mathbb{Q}/\mathbb{Z} = \mathbb{Q}/\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Z}$ this is the statement that multiplication by 2 has the element $1/2$ in its kernel).
- (4) The direct sum of flat modules is flat (Exercise 5). In particular, $\mathbb{Q} \oplus \mathbb{Z}$ is flat. This module is neither projective nor injective (since \mathbb{Q} is not projective by Exercise 8 and \mathbb{Z} is not injective by Proposition 36 (cf. Exercises 3 and 4)).

We close this section with an important relation between Hom and tensor products:

Theorem 43. (Adjoint Associativity) Let R and S be rings, let A be a right R -module, let B be an (R, S) -bimodule and let C be a right S -module. Then there is an isomorphism of abelian groups:

$$\text{Hom}_S(A \otimes_R B, C) \cong \text{Hom}_R(A, \text{Hom}_S(B, C))$$

(the homomorphism groups are right module homomorphisms—note that $\text{Hom}_S(B, C)$ has the structure of a right R -module, cf. the exercises). If $R = S$ is commutative this is an isomorphism of R -modules with the standard R -module structures.

Proof: Suppose $\varphi : A \otimes_R B \rightarrow C$ is a homomorphism. For any fixed $a \in A$ define the map $\Phi(a)$ from B to C by $\Phi(a)(b) = \varphi(a \otimes b)$. It is easy to check that $\Phi(a)$ is a homomorphism of right S -modules and that the map Φ from A to $\text{Hom}_S(B, C)$ given by mapping a to $\Phi(a)$ is a homomorphism of right R -modules. Then $f(\varphi) = \Phi$ defines a group homomorphism from $\text{Hom}_S(A \otimes_R B, C)$ to $\text{Hom}_R(A, \text{Hom}_S(B, C))$. Conversely, suppose $\Phi : A \rightarrow \text{Hom}_S(B, C)$ is a homomorphism. The map from $A \times B$ to C defined by mapping (a, b) to $\Phi(a)(c)$ is an R -balanced map, so induces a homomorphism φ from $A \otimes_R B$ to C . Then $g(\Phi) = \varphi$ defines a group homomorphism inverse to f and gives the isomorphism in the theorem.

As a first application of Theorem 43 we give an alternate proof of the first result in Theorem 39 that the tensor product is right exact in the case where $S = R$ is a commutative ring. If $0 \rightarrow L \rightarrow M \rightarrow N \rightarrow 0$ is exact, then by Theorem 33 the sequence

$$0 \rightarrow \text{Hom}_R(N, E) \rightarrow \text{Hom}_R(M, E) \rightarrow \text{Hom}_R(L, E)$$

is exact for every R -module E . Then by Theorem 28, the sequence

$$0 \rightarrow \text{Hom}_R(D, \text{Hom}_R(N, E)) \rightarrow \text{Hom}_R(D, \text{Hom}_R(M, E)) \rightarrow \text{Hom}_R(D, \text{Hom}_R(L, E))$$

is exact for all D and all E . By adjoint associativity, this means the sequence

$$0 \rightarrow \text{Hom}_R(D \otimes_R N, E) \rightarrow \text{Hom}_R(D \otimes_R M, E) \rightarrow \text{Hom}_R(D \otimes_R L, E)$$

is exact for any D and all E . Then, by the second part of Theorem 33, it follows that the sequence

$$D \otimes_R L \rightarrow D \otimes_R M \rightarrow D \otimes_R N \rightarrow 0$$

is exact for all D , which is the right exactness of the tensor product.

As a second application of Theorem 43 we prove that the tensor product of two projective modules over a commutative ring R is again projective (see also Exercise 9 for a more direct proof).

Corollary 44. If R is commutative then the tensor product of two projective R -modules is projective.

Proof: Let P_1 and P_2 be projective modules. Then by Corollary 32, $\text{Hom}_R(P_2, \underline{\quad})$ is an exact functor from the category of R -modules to the category of R -modules. Then the composition $\text{Hom}_R(P_1, \text{Hom}_R(P_2, \underline{\quad}))$ is an exact functor by the same corollary. By Theorem 43 this means that $\text{Hom}_R(P_1 \otimes_R P_2, \underline{\quad})$ is an exact functor on R -modules. It follows again from Corollary 32 that $P_1 \otimes_R P_2$ is projective.

Summary

Each of the functors $\text{Hom}_R(A, \underline{\quad})$, $\text{Hom}_R(\underline{\quad}, A)$, and $A \otimes_R \underline{\quad}$, map left R -modules to abelian groups; the functor $\underline{\quad} \otimes_R A$ maps right R -modules to abelian groups. When R is commutative all four functors map R -modules to R -modules.

- (1) Let A be a left R -module. The functor $\text{Hom}_R(A, \underline{\quad})$ is covariant and left exact; the module A is projective if and only if $\text{Hom}_R(A, \underline{\quad})$ is exact (i.e., is also right exact).

- (2) Let A be a left R -module. The functor $\text{Hom}_R(_, A)$ is contravariant and left exact; the module A is injective if and only if $\text{Hom}_R(_, A)$ is exact.
- (3) Let A be a right R -module. The functor $A \otimes_R _$ is covariant and right exact; the module A is flat if and only if $A \otimes_R _$ is exact (i.e., is also left exact).
- (4) Let A be a left R -module. The functor $_ \otimes_R A$ is covariant and right exact; the module A is flat if and only if $_ \otimes_R A$ is exact.
- (5) Projective modules are flat. The \mathbb{Z} -module \mathbb{Q}/\mathbb{Z} is injective but not flat. The \mathbb{Z} -module $\mathbb{Z} \oplus \mathbb{Q}$ is flat but neither projective nor injective.

EXERCISES

Let R be a ring with 1.

1. Suppose that

$$\begin{array}{ccccc} A & \xrightarrow{\psi} & B & \xrightarrow{\varphi} & C \\ \alpha \downarrow & & \beta \downarrow & & \gamma \downarrow \\ A' & \xrightarrow{\psi'} & B' & \xrightarrow{\varphi'} & C' \end{array}$$

is a commutative diagram of groups and that the rows are exact. Prove that

- (a) if φ and α are surjective, and β is injective then γ is injective. [If $c \in \ker \gamma$, show there is a $b \in B$ with $\varphi(b) = c$. Show that $\varphi'(\beta(b)) = 0$ and deduce that $\beta(b) = \psi'(a')$ for some $a' \in A'$. Show there is an $a \in A$ with $\alpha(a) = a'$ and that $\beta(\psi(a)) = \beta(b)$. Conclude that $b = \psi(a)$ and hence $c = \varphi(b) = 0$.]
- (b) if ψ' , α , and γ are injective, then β is injective,
- (c) if φ , α , and γ are surjective, then β is surjective,
- (d) if β is injective, α and γ are surjective, then γ is injective,
- (e) if β is surjective, γ and ψ' are injective, then α is surjective.

2. Suppose that

$$\begin{array}{ccccccc} A & \longrightarrow & B & \longrightarrow & C & \longrightarrow & D \\ \alpha \downarrow & & \beta \downarrow & & \gamma \downarrow & & \delta \downarrow \\ A' & \longrightarrow & B' & \longrightarrow & C' & \longrightarrow & D' \end{array}$$

is a commutative diagram of groups, and that the rows are exact. Prove that

- (a) if α is surjective, and β , δ are injective, then γ is injective.
- (b) if δ is injective, and α , γ are surjective, then β is surjective.

- 3. Let P_1 and P_2 be R -modules. Prove that $P_1 \oplus P_2$ is a projective R -module if and only if both P_1 and P_2 are projective.
- 4. Let Q_1 and Q_2 be R -modules. Prove that $Q_1 \oplus Q_2$ is an injective R -module if and only if both Q_1 and Q_2 are injective.
- 5. Let A_1 and A_2 be R -modules. Prove that $A_1 \oplus A_2$ is a flat R -module if and only if both A_1 and A_2 are flat. More generally, prove that an arbitrary direct sum $\sum A_i$ of R -modules is flat if and only if each A_i is flat. [Use the fact that tensor product commutes with arbitrary direct sums.]
- 6. Prove that the following are equivalent for a ring R :
 - (i) Every R -module is projective.
 - (ii) Every R -module is injective.

7. Let A be a nonzero finite abelian group.
- Prove that A is not a projective \mathbb{Z} -module.
 - Prove that A is not an injective \mathbb{Z} -module.
8. Let Q be a nonzero divisible \mathbb{Z} -module. Prove that Q is not a projective \mathbb{Z} -module. Deduce that the rational numbers \mathbb{Q} is not a projective \mathbb{Z} -module. [Show first that if F is any free module then $\bigcap_{n=1}^{\infty} nF = 0$ (use a basis of F to prove this). Now suppose to the contrary that Q is projective and derive a contradiction from Proposition 30(4).]
9. Assume R is commutative with 1.
- Prove that the tensor product of two free R -modules is free. [Use the fact that tensor products commute with direct sums.]
 - Use (a) to prove that the tensor product of two projective R -modules is projective.
10. Let R and S be rings with 1 and let M and N be left R -modules. Assume also that M is an (R, S) -bimodule.
- For $s \in S$ and for $\varphi \in \text{Hom}_R(M, N)$ define $(s\varphi) : M \rightarrow N$ by $(s\varphi)(m) = \varphi(ms)$. Prove that $s\varphi$ is a homomorphism of left R -modules, and that this action of S on $\text{Hom}_R(M, N)$ makes it into a *left* S -module.
 - Let $S = R$ and let $M = R$ (considered as an (R, R) -bimodule by left and right ring multiplication on itself). For each $n \in N$ define $\varphi_n : R \rightarrow N$ by $\varphi_n(r) = rn$, i.e., φ_n is the unique R -module homomorphism mapping 1_R to n . Show that $\varphi_n \in \text{Hom}_R(R, N)$. Use part (a) to show that the map $n \mapsto \varphi_n$ is an isomorphism of left R -modules: $N \cong \text{Hom}_R(R, N)$.
 - Deduce that if N is a free (respectively, projective, injective, flat) left R -module, then $\text{Hom}_R(R, N)$ is also a free (respectively, projective, injective, flat) left R -module.
11. Let R and S be rings with 1 and let M and N be left R -modules. Assume also that N is an (R, S) -bimodule.
- For $s \in S$ and for $\varphi \in \text{Hom}_R(M, N)$ define $(\varphi s) : M \rightarrow N$ by $(\varphi s)(m) = \varphi(m)s$. Prove that φs is a homomorphism of left R -modules, and that this action of S on $\text{Hom}_R(M, N)$ makes it into a *right* S -module. Deduce that $\text{Hom}_R(M, R)$ is a right R -module, for any R -module M —called the *dual module* to M .
 - Let $N = R$ be considered as an (R, R) -bimodule as usual. Under the action defined in part (a) show that the map $r \mapsto \varphi_r$ is an isomorphism of right R -modules: $\text{Hom}_R(R, R) \cong R$, where φ_r is the homomorphism that maps 1_R to r . Deduce that if M is a finitely generated free left R -module, then $\text{Hom}_R(M, R)$ is a free right R -module of the same rank. (cf. also Exercise 13.)
 - Show that if M is a finitely generated projective R -module then its dual module $\text{Hom}_R(M, R)$ is also projective.
12. Let A be an R -module, let I be any nonempty index set and for each $i \in I$ let B_i be an R -module. Prove the following isomorphisms of abelian groups; when R is commutative prove also that these are R -module isomorphisms. (Arbitrary direct sums and direct products of modules are introduced in Exercise 20 of Section 3.)
- $\text{Hom}_R(\bigoplus_{i \in I} B_i, A) \cong \prod_{i \in I} \text{Hom}_R(B_i, A)$
 - $\text{Hom}_R(A, \prod_{i \in I} B_i) \cong \prod_{i \in I} \text{Hom}_R(A, B_i)$.
13. (a) Show that the dual of the free \mathbb{Z} -module with countable basis is not free. [Use the preceding exercise and Exercise 24, Section 3.] (See also Exercise 5 in Section 11.3.)
(b) Show that the dual of the free \mathbb{Z} -module with countable basis is also not projective. [You may use the fact that any submodule of a free \mathbb{Z} -module is free.]
14. Let $0 \longrightarrow L \xrightarrow{\psi} M \xrightarrow{\varphi} N \longrightarrow 0$ be a sequence of R -modules.

- (a) Prove that the associated sequence

$$0 \longrightarrow \text{Hom}_R(D, L) \xrightarrow{\psi'} \text{Hom}_R(D, M) \xrightarrow{\varphi'} \text{Hom}_R(D, N) \longrightarrow 0$$

is a short exact sequence of abelian groups for all R -modules D if and only if the original sequence is a split short exact sequence. [To show the sequence splits, take $D = N$ and show the lift of the identity map in $\text{Hom}_R(N, N)$ to $\text{Hom}_R(N, M)$ is a splitting homomorphism for φ .]

- (b) Prove that the associated sequence

$$0 \longrightarrow \text{Hom}_R(N, D) \xrightarrow{\varphi'} \text{Hom}_R(M, D) \xrightarrow{\psi'} \text{Hom}_R(L, D) \longrightarrow 0$$

is a short exact sequence of abelian groups for all R -modules D if and only if the original sequence is a split short exact sequence.

15. Let M be a left R -module where R is a ring with 1.

- (a) Show that $\text{Hom}_{\mathbb{Z}}(R, M)$ is a left R -module under the action $(r\varphi)(r') = \varphi(r'r)$ (see Exercise 10).
- (b) Suppose that $0 \rightarrow A \xrightarrow{\psi} B$ is an exact sequence of R -modules. Prove that if every homomorphism f from A to M lifts to a homomorphism F from B to M with $f = F \circ \psi$, then every homomorphism f' from A to $\text{Hom}_{\mathbb{Z}}(R, M)$ lifts to a homomorphism F' from B to $\text{Hom}_{\mathbb{Z}}(R, M)$ with $f' = F' \circ \psi$. [Given f' , show that $f(a) = f'(a)(1_R)$ defines a homomorphism of A to M . If F is the associated lift of f to B , show that $F'(b)(r) = F(rb)$ defines a homomorphism from B to $\text{Hom}_{\mathbb{Z}}(R, M)$ that lifts f' .]
- (c) Prove that if Q is an injective R -module then $\text{Hom}_{\mathbb{Z}}(R, Q)$ is also an injective R -module.

16. This exercise proves Theorem 38 that every left R -module M is contained in an injective left R -module.

- (a) Show that M is contained in an injective \mathbb{Z} -module Q . [M is a \mathbb{Z} -module—use Corollary 37.]
- (b) Show that $\text{Hom}_R(R, M) \subseteq \text{Hom}_{\mathbb{Z}}(R, M) \subseteq \text{Hom}_{\mathbb{Z}}(R, Q)$.
- (c) Use the R -module isomorphism $M \cong \text{Hom}_R(R, M)$ (Exercise 10) and the previous exercise to conclude that M is contained in an injective module.

17. This exercise completes the proof of Proposition 34. Suppose that Q is an R -module with the property that every short exact sequence $0 \rightarrow Q \rightarrow M_1 \rightarrow N \rightarrow 0$ splits and suppose that the sequence $0 \rightarrow L \xrightarrow{\psi} M$ is exact. Prove that every R -module homomorphism f from L to Q can be lifted to an R -module homomorphism F from M to Q with $f = F \circ \psi$. [By the previous exercise, Q is contained in an injective R -module. Use the splitting property together with Exercise 4 (noting that Exercise 4 can be proved using (2) in Proposition 34 as the definition of an injective module).]

18. Prove that the injective hull of the \mathbb{Z} -module \mathbb{Z} is \mathbb{Q} . [Let H be the injective hull of \mathbb{Z} and argue that \mathbb{Q} contains an isomorphic copy of H . Use the divisibility of H to show $1/n \in H$ for all nonzero integers n , and deduce that $H = \mathbb{Q}$.]

19. If F is a field, prove that the injective hull of F is F .

20. Prove that the polynomial ring $R[x]$ in the indeterminate x over the commutative ring R is a flat R -module.

21. Let R and S be rings with 1 and suppose M is a right R -module, and N is an (R, S) -bimodule. If M is flat over R and N is flat as an S -module prove that $M \otimes_R N$ is flat as a right S -module.

22. Suppose that R is a commutative ring and that M and N are flat R -modules. Prove that $M \otimes_R N$ is a flat R -module. [Use the previous exercise.]
23. Prove that the (right) module $M \otimes_R S$ obtained by changing the base from the ring R to the ring S (by some homomorphism $f : R \rightarrow S$ with $f(1_R) = 1_S$, cf. Example 6 following Corollary 12 in Section 4) of the flat (right) R -module M is a flat S -module.
24. Prove that A is a flat R -module if and only if for any left R -modules L and M where L is *finitely generated*, then $\psi : L \rightarrow M$ injective implies that also $1 \otimes \psi : A \otimes_R L \rightarrow A \otimes_R M$ is injective. [Use the techniques in the proof of Corollary 42.]
25. (A Flatness Criterion) Parts (a)-(c) of this exercise prove that A is a flat R -module if and only if for every finitely generated ideal I of R , the map from $A \otimes_R I \rightarrow A \otimes_R R \cong A$ induced by the inclusion $I \subseteq R$ is again injective (or, equivalently, $A \otimes_R I \cong AI \subseteq A$).
- Prove that if A is flat then $A \otimes_R I \rightarrow A \otimes_R R$ is injective.
 - If $A \otimes_R I \rightarrow A \otimes_R R$ is injective for every finitely generated ideal I , prove that $A \otimes_R I \rightarrow A \otimes_R R$ is injective for every ideal I . Show that if K is any submodule of a finitely generated free module F then $A \otimes_R K \rightarrow A \otimes_R F$ is injective. Show that the same is true for any free module F . [Cf. the proof of Corollary 42.]
 - Under the assumption in (b), suppose L and M are R -modules and $L \xrightarrow{\psi} M$ is injective. Prove that $A \otimes_R L \xrightarrow{1 \otimes \psi} A \otimes_R M$ is injective and conclude that A is flat. [Write M as a quotient of the free module F , giving a short exact sequence
- $$0 \longrightarrow K \longrightarrow F \xrightarrow{f} M \longrightarrow 0.$$
- Show that if $J = f^{-1}(\psi(L))$ and $\iota : J \rightarrow F$ is the natural injection, then the diagram
- $$\begin{array}{ccccccc} 0 & \longrightarrow & K & \longrightarrow & J & \longrightarrow & L & \longrightarrow & 0 \\ & & id \downarrow & & \iota \downarrow & & \psi \downarrow & & \\ 0 & \longrightarrow & K & \longrightarrow & F & \longrightarrow & M & \longrightarrow & 0 \end{array}$$
- is commutative with exact rows. Show that the induced diagram
- $$\begin{array}{ccccccc} A \otimes_R K & \longrightarrow & A \otimes_R J & \longrightarrow & A \otimes_R L & \longrightarrow & 0 \\ id \downarrow & & 1 \otimes \iota \downarrow & & 1 \otimes \psi \downarrow & & \\ A \otimes_R K & \longrightarrow & A \otimes_R F & \longrightarrow & A \otimes_R M & \longrightarrow & 0 \end{array}$$
- is commutative with exact rows. Use (b) to show that $1 \otimes \iota$ is injective, then use Exercise 1 to conclude that $1 \otimes \psi$ is injective.]
- (A Flatness Criterion for quotients) Suppose $A = F/K$ where F is flat (e.g., if F is free) and K is an R -submodule of F . Prove that A is flat if and only if $FI \cap K = KI$ for every finitely generated ideal I of R . [Use (a) to prove $F \otimes_R I \cong FI$ and observe the image of $K \otimes_R I$ is KI ; tensor the exact sequence $0 \rightarrow K \rightarrow F \rightarrow A \rightarrow 0$ with I to prove that $A \otimes_R I \cong FI/KI$, and apply the flatness criterion.]
26. Suppose R is a P.I.D. This exercise proves that A is a flat R -module if and only if A is torsion free R -module (i.e., if $a \in A$ is nonzero and $r \in R$, then $ra = 0$ implies $r = 0$).
- Suppose that A is flat and for fixed $r \in R$ consider the map $\psi_r : R \rightarrow R$ defined by multiplication by r : $\psi_r(x) = rx$. If r is nonzero show that ψ_r is an injection. Conclude from the flatness of A that the map from A to A defined by mapping a to ra is injective and that A is torsion free.
 - Suppose that A is torsion free. If I is a nonzero ideal of R , then $I = rR$ for some nonzero $r \in R$. Show that the map ψ_r in (a) induces an isomorphism $R \cong I$ of

R -modules and that the composite $R \xrightarrow{\psi} I \xrightarrow{\iota} R$ of ψ_r with the inclusion $\iota : I \subseteq R$ is multiplication by r . Prove that the composite $A \otimes_R R \xrightarrow{1 \otimes \psi_r} A \otimes_R I \xrightarrow{1 \otimes \iota} A \otimes_R R$ corresponds to the map $a \mapsto ra$ under the identification $A \otimes_R R = A$ and that this composite is injective since A is torsion free. Show that $1 \otimes \psi_r$ is an isomorphism and deduce that $1 \otimes \iota$ is injective. Use the previous exercise to conclude that A is flat.

27. Let M , A and B be R -modules.

- (a) Suppose $f : A \rightarrow M$ and $g : B \rightarrow M$ are R -module homomorphisms. Prove that $X = \{(a, b) \mid a \in A, b \in B \text{ with } f(a) = g(b)\}$ is an R -submodule of the direct sum $A \oplus B$ (called the *pullback* or *fiber product* of f and g) and that there is a commutative diagram

$$\begin{array}{ccc} X & \xrightarrow{\pi_2} & B \\ \pi_1 \downarrow & & \downarrow g \\ A & \xrightarrow{f} & M \end{array}$$

where π_1 and π_2 are the natural projections onto the first and second components.

- (b) Suppose $f' : M \rightarrow A$ and $g' : M \rightarrow B$ are R -module homomorphisms. Prove that the quotient Y of $A \oplus B$ by $\{(f'(m), -g'(m)) \mid m \in M\}$ is an R -module (called the *pushout* or *fiber sum* of f' and g') and that there is a commutative diagram

$$\begin{array}{ccc} M & \xrightarrow{g'} & B \\ f' \downarrow & & \downarrow \pi'_2 \\ A & \xrightarrow{\pi'_1} & Y \end{array}$$

where π'_1 and π'_2 are the natural maps to the quotient induced by the maps into the first and second components.

28. (a) (*Schanuel's Lemma*) If $0 \rightarrow K \rightarrow P \xrightarrow{\varphi} M \rightarrow 0$ and $0 \rightarrow K' \rightarrow P' \xrightarrow{\varphi'} M \rightarrow 0$ are exact sequences of R -modules where P and P' are projective, prove $P \oplus K' \cong P' \oplus K$ as R -modules. [Show that there is an exact sequence $0 \rightarrow \ker \pi \rightarrow X \xrightarrow{\pi} P \rightarrow 0$ with $\ker \pi \cong K'$, where X is the fiber product of φ and φ' as in the previous exercise. Deduce that $X \cong P \oplus K'$. Show similarly that $X \cong P' \oplus K$.]
- (b) If $0 \rightarrow M \rightarrow Q \xrightarrow{\psi} L \rightarrow 0$ and $0 \rightarrow M \rightarrow Q' \xrightarrow{\psi'} L' \rightarrow 0$ are exact sequences of R -modules where Q and Q' are injective, prove $Q \oplus L' \cong Q' \oplus L$ as R -modules.

The R -modules M and N are said to be *projectively equivalent* if $M \oplus P \cong N \oplus P'$ for some projective modules P , P' . Similarly, M and N are *injectively equivalent* if $M \oplus Q \cong N \oplus Q'$ for some injective modules Q , Q' . The previous exercise shows K and K' are projectively equivalent and L and L' are injectively equivalent.

CHAPTER 11

Vector Spaces

In this chapter we review the basic theory of finite dimensional vector spaces over an arbitrary field F (some infinite dimensional vector space theory is covered in the exercises). Since the proofs are identical to the corresponding arguments for real vector spaces our treatment is very terse. For the most part we include only those results which are used in other parts of the text so basic topics such as Gauss–Jordan elimination, row echelon forms, methods for finding bases of subspaces, elementary properties of matrices, etc., are not covered or are discussed in the exercises. The reader should therefore consider this chapter as a refresher in linear algebra and as a prelude to field theory and Galois theory. Characteristic polynomials and eigenvalues will be reviewed and treated in a larger context in the next chapter.

11.1 DEFINITIONS AND BASIC THEORY

The terminology for vector spaces is slightly different from that of modules, that is, when the ring R is a field there are different names for many of the properties of R -modules which we defined in the last chapter. The following is a dictionary of these new terms (many of which may already be familiar). The definition of each corresponding vector space property is the same (verbatim) as the module-theoretic definition with the only added assumption being that the ring R is a field (so these definitions are not repeated here).

Terminology for R any Ring

- M is an R -module
- m is an element of M
- α is a ring element
- N is a submodule of M
- M/N is a quotient module
- M is a free module of rank n
- M is a finitely generated module
- M is a nonzero cyclic module
- $\varphi : M \rightarrow N$ is an R -module homomorphism
- M and N are isomorphic as R -modules
- the subset A of M generates M
- $M = RA$

Terminology for R a Field

- M is a vector space over R
- m is a vector in M
- α is a scalar
- N is a subspace of M
- M/N is a quotient space
- M is a vector space of dimension n
- M is a finite dimensional vector space
- M is a 1-dimensional vector space
- $\varphi : M \rightarrow N$ is a linear transformation
- M and N are isomorphic vector spaces
- the subset A of M spans M
- each element of M is a linear combination of elements of A i.e., $M = \text{Span}(A)$

For the remainder of this chapter F is a field and V is a vector space over F .

One of the first results we shall prove about vector spaces is that they are free F -modules, that is, they have bases. Although our arguments treat only the case of finite dimensional spaces, the corresponding result for arbitrary vector spaces is proved in the exercises as an application of Zorn's Lemma. The reader may first wish to review the section in the previous chapter on free modules, especially their properties pertaining to homomorphisms.

Definition.

- (1) A subset S of V is called a set of *linearly independent* vectors if an equation $\alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_n v_n = 0$ with $\alpha_1, \alpha_2, \dots, \alpha_n \in F$ and $v_1, v_2, \dots, v_n \in S$ implies $\alpha_1 = \alpha_2 = \cdots = \alpha_n = 0$.
- (2) A *basis* of a vector space V is an ordered set of linearly independent vectors which span V . In particular two bases will be considered different even if one is simply a rearrangement of the other. This is sometimes referred to as an *ordered basis*.

Examples

- (1) The space $V = F[x]$ of polynomials in the variable x with coefficients from the field F is in particular a vector space over F . The elements $1, x, x^2, \dots$ are linearly independent by definition (i.e., a polynomial is 0 if and only if all its coefficients are 0). Since these elements also span V by definition, they are a basis for V .
- (2) The collection of solutions of a linear, homogeneous, constant coefficient differential equation (for example, $y'' - 3y' + 2y = 0$) over \mathbb{C} form a vector space over \mathbb{C} since differentiation is a linear operator. Elements of this vector space are linearly independent if they are linearly independent as functions. For example, e^t and e^{2t} are easily seen to be solutions of the equation $y'' - 3y' + 2y = 0$ (differentiation with respect to t). They are linearly independent functions since $ae^t + be^{2t} = 0$ implies $a + b = 0$ (let $t = 0$) and $ae + be^2 = 0$ (let $t = 1$) and the only solution to these two equations is $a = b = 0$. It is a theorem in differential equations that these elements span the set of solutions of this equation, hence are a basis for this space.

Proposition 1. Assume the set $\mathcal{A} = \{v_1, v_2, \dots, v_n\}$ spans the vector space V but no proper subset of \mathcal{A} spans V . Then \mathcal{A} is a basis of V . In particular, any finitely generated (i.e., finitely spanned) vector space over F is a free F -module.

Proof: It is only necessary to prove that v_1, v_2, \dots, v_n are linearly independent. Suppose $\alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_n v_n = 0$ where not all of the α_i are 0. By reordering, we may assume that $\alpha_1 \neq 0$ and then

$$v_1 = -\frac{1}{\alpha_1}(\alpha_2 v_2 + \cdots + \alpha_n v_n).$$

It follows that $\{v_2, v_3, \dots, v_n\}$ also spans V since any linear combination of v_1, v_2, \dots, v_n can be written as a linear combination of v_2, v_3, \dots, v_n using the equation above. This is a contradiction.

Example

Let F be a field and consider $F[x]/(f(x))$ where $f(x) = x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0$. The ideal $(f(x))$ is a subspace of the vector space $F[x]$ and the quotient $F[x]/(f(x))$ is also a vector space over F . By the Euclidean Algorithm, every polynomial $a(x) \in F[x]$ can be written uniquely in the form $a(x) = q(x)f(x) + r(x)$ where $r(x) \in F[x]$ and $0 \leq \deg r(x) \leq n - 1$. Since $q(x)f(x) \in (f(x))$, it follows that every element of the quotient is represented by a polynomial $r(x)$ of degree $\leq n - 1$. Two distinct such polynomials cannot be the same in the quotient since this would say their difference (which is a nonzero polynomial of degree at most $n - 1$) would be divisible by $f(x)$ (which is of degree n). It follows that the elements $\bar{1}, \bar{x}, \bar{x^2}, \dots, \bar{x^{n-1}}$ (the bar denotes the image of these elements in the quotient, as usual) span $F[x]/(f(x))$ as a vector space over F and that no proper subset of these elements also spans, hence these elements give a basis for $F[x]/(f(x))$.

Corollary 2. Assume the finite set \mathcal{A} spans the vector space V . Then \mathcal{A} contains a basis of V .

Proof: Any subset \mathcal{B} of \mathcal{A} spanning V such that no proper subset of \mathcal{B} also spans V (there clearly exist such subsets) is a basis for V by Proposition 1.

Theorem 3. (A Replacement Theorem) Assume $\mathcal{A} = \{a_1, a_2, \dots, a_n\}$ is a basis for V containing n elements and $\{b_1, b_2, \dots, b_m\}$ is a set of linearly independent vectors in V . Then there is an ordering a_1, a_2, \dots, a_n such that for each $k \in \{1, 2, \dots, m\}$ the set $\{b_1, b_2, \dots, b_k, a_{k+1}, a_{k+2}, \dots, a_n\}$ is a basis of V . In other words, the elements b_1, b_2, \dots, b_m can be used to successively replace the elements of the basis \mathcal{A} , still retaining a basis. In particular, $n \geq m$.

Proof: Proceed by induction on k . If $k = 0$ there is nothing to prove, since \mathcal{A} is given as a basis for V . Suppose now that $\{b_1, b_2, \dots, b_k, a_{k+1}, a_{k+2}, \dots, a_n\}$ is a basis for V . Then in particular this is a spanning set, so b_{k+1} is a linear combination:

$$b_{k+1} = \beta_1 b_1 + \dots + \beta_k b_k + \alpha_{k+1} a_{k+1} + \dots + \alpha_n a_n. \quad (11.1)$$

Not all of the α_i can be 0, since this would imply b_{k+1} is a linear combination of b_1, b_2, \dots, b_k , contrary to the linear independence of these elements. By reordering if necessary, we may assume $\alpha_{k+1} \neq 0$. Then solving this last equation for a_{k+1} as a linear combination of b_{k+1} and $b_1, b_2, \dots, b_k, a_{k+2}, \dots, a_n$ shows

$$\text{Span}\{b_1, b_2, \dots, b_k, b_{k+1}, a_{k+2}, \dots, a_n\} = \text{Span}\{b_1, b_2, \dots, b_k, a_{k+1}, a_{k+2}, \dots, a_n\}$$

and so this is a spanning set for V . It remains to show $b_1, \dots, b_k, b_{k+1}, a_{k+2}, \dots, a_n$ are linearly independent. If

$$\beta_1 b_1 + \dots + \beta_k b_k + \beta_{k+1} b_{k+1} + \alpha_{k+2} a_{k+2} + \dots + \alpha_n a_n = 0 \quad (11.2)$$

then substituting for b_{k+1} from the expression for b_{k+1} in equation (1), we obtain a linear combination of $\{b_1, b_2, \dots, b_k, a_{k+1}, a_{k+2}, \dots, a_n\}$ equal to 0, where the coefficient of a_{k+1} is β_{k+1} . Since this last set is a basis by induction, all the coefficients in this linear combination, in particular β_{k+1} , must be 0. But then equation (2) is

$$\beta_1 b_1 + \dots + \beta_k b_k + \alpha_{k+2} a_{k+2} + \dots + \alpha_n a_n = 0.$$

Again by the induction hypothesis all the other coefficients must be 0 as well. Thus $\{b_1, b_2, \dots, b_k, b_{k+1}, a_{k+2}, \dots, a_n\}$ is a basis for V , and the induction is complete.

Corollary 4.

- (1) Suppose V has a finite basis with n elements. Any set of linearly independent vectors has $\leq n$ elements. Any spanning set has $\geq n$ elements.
- (2) If V has some finite basis then any two bases of V have the same cardinality.

Proof: (1) This is a restatement of the last result of Theorem 3 and Corollary 2.
(2) This is immediate from (1) since a basis is both a spanning set and a linearly independent set.

Definition. If V is a finitely generated F -module (i.e., has a finite basis) the cardinality of any basis is called the *dimension* of V and is denoted by $\dim_F V$, or just $\dim V$ when F is clear from the context, and V is said to be *finite dimensional* over F . If V is not finitely generated, V is said to be infinite dimensional (written $\dim V = \infty$).

Examples

- (1) The dimension of the space of solutions to the differential equation $y'' - 3y' + 2y = 0$ over \mathbb{C} is 2 (with basis e^t, e^{2t} , for example). In general, it is a theorem in differential equations that the space of solutions of an n^{th} order linear, homogeneous, constant coefficient differential equation of degree n over \mathbb{C} form a vector space over \mathbb{C} of dimension n .
- (2) The dimension over F of the quotient $F[x]/(f(x))$ by the nonzero polynomial $f(x)$ considered above is $n = \deg f(x)$. The space $F[x]$ and its subspace $(f(x))$ are infinite dimensional vector spaces over F .

Corollary 5. (Building–Up Lemma) If A is a set of linearly independent vectors in the finite dimensional space V then there exists a basis of V containing A .

Proof: This is also immediate from Theorem 3, since we can use the elements of A to successively replace the elements of any given basis for V (which exists by the assumption that V is finite dimensional).

Theorem 6. If V is an n dimensional vector space over F , then $V \cong F^n$. In particular, any two finite dimensional vector spaces over F of the same dimension are isomorphic.

Proof: Let v_1, v_2, \dots, v_n be a basis for V . Define the map

$$\varphi : F^n \rightarrow V \quad \text{by} \quad \varphi(\alpha_1, \alpha_2, \dots, \alpha_n) = \alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_n v_n.$$

The map φ is clearly F -linear, is surjective since the v_i span V , and is injective since the v_i are linearly independent, hence is an isomorphism.

Examples

- (1) Let \mathbb{F} be a finite field with q elements and let W be a k -dimensional vector space over \mathbb{F} . We show that the number of distinct bases of W is

$$(q^k - 1)(q^k - q)(q^k - q^2) \dots (q^k - q^{k-1}).$$

Every basis of W can be built up as follows. Any nonzero vector w_1 can be the first element of a basis. Since W is isomorphic to \mathbb{F}^k , $|W| = q^k$, so there are $q^k - 1$ choices for w_1 . Any vector not in the 1-dimensional space spanned by w_1 is linearly independent from w_1 and so may be chosen for the second basis element, w_2 . A 1-dimensional space is isomorphic to \mathbb{F} and so has q elements. Thus there are $q^k - q$ choices for w_2 . Proceeding in this way one sees that at the i^{th} stage any vector not in the $(i-1)$ -dimensional space spanned by w_1, w_2, \dots, w_{i-1} will be linearly independent from w_1, w_2, \dots, w_{i-1} and so may be chosen for the i^{th} basis vector w_i . An $(i-1)$ -dimensional space is isomorphic to \mathbb{F}^{i-1} and so has q^{i-1} elements. Thus there are $q^k - q^{i-1}$ choices for w_i . The process terminates when w_k is chosen, for then we have k linear independent vectors in a k -dimensional space, hence a basis.

- (2) Let \mathbb{F} be a finite field with q elements and let V be an n -dimensional vector space over \mathbb{F} . For each $k \in \{1, 2, \dots, n\}$ we show that the number of subspaces of V of dimension k is

$$\frac{(q^n - 1)(q^n - q) \dots (q^n - q^{k-1})}{(q^k - 1)(q^k - q) \dots (q^k - q^{k-1})}.$$

Any k -dimensional space is spanned by k independent vectors. By arguing as in the preceding example the numerator of the above expression is the number of ways of picking k independent vectors from an n -dimensional space. Two sets of k independent vectors span the same space W if and only if they are both bases of the k -dimensional space W . In order to obtain the formula for the number of distinct subspaces of dimension k we must divide by the number of repetitions, i.e., the number of bases of a fixed k -dimensional space. This factor which appears in the denominator is precisely the number computed in Example 1.

Next, we prove an important relation between the dimension of a subspace, the dimension of its associated quotient space and the dimension of the whole space:

Theorem 7. Let V be a vector space over F and let W be a subspace of V . Then V/W is a vector space with $\dim V = \dim W + \dim V/W$ (where if one side is infinite then both are).

Proof: Suppose W has dimension m and V has dimension n over F and let w_1, w_2, \dots, w_m be a basis for W . By Corollary 5, these linearly independent elements of V can be extended to a basis $w_1, w_2, \dots, w_m, v_{m+1}, \dots, v_n$ of V . The natural surjective projection map of V into V/W maps each w_i to 0. No linear combination of the v_i is mapped to 0, since this would imply this linear combination is an element of W , contrary to the choice of the v_i . Hence, the image V/W of this projection map is isomorphic to the subspace of V spanned by the v_i , hence $\dim V/W = n - m$, which is the theorem when the dimensions are finite. If either side is infinite it is an easy exercise to produce an infinite number of linearly independent vectors showing the other side is also infinite.

Corollary 8. Let $\varphi : V \rightarrow U$ be a linear transformation of vector spaces over F . Then $\ker \varphi$ is a subspace of V , $\varphi(V)$ is a subspace of U and $\dim V = \dim \ker \varphi + \dim \varphi(V)$.

Proof: This follows immediately from Theorem 7. Note that the proof of Theorem 7 is in fact the special case of Corollary 8 where U is the quotient V/W and φ is the natural projection homomorphism.

Corollary 9. Let $\varphi : V \rightarrow W$ be a linear transformation of vector spaces of the same finite dimension. Then the following are equivalent:

- (1) φ is an isomorphism
- (2) φ is injective, i.e., $\ker \varphi = 0$
- (3) φ is surjective, i.e., $\varphi(V) = W$
- (4) φ sends a basis of V to a basis of W .

Proof: The equivalence of these conditions follows from Corollary 8 by counting dimensions.

Definition. If $\varphi : V \rightarrow U$ is a linear transformation of vector spaces over F , $\ker \varphi$ is sometimes called the *null space* of φ and the dimension of $\ker \varphi$ is called the *nullity* of φ . The dimension of $\varphi(V)$ is called the *rank* of φ . If $\ker \varphi = 0$, the transformation is said to be *nonsingular*.

Example

Let F be a finite field with q elements and let V be an n -dimensional vector space over F . Recall that the *general linear group* $GL(V)$ is the group of all nonsingular linear transformations from V to V (the group operation being composition). We show that the order of this group is

$$|GL(V)| = (q^n - 1)(q^n - q)(q^n - q^2) \dots (q^n - q^{n-1}).$$

To see this, fix a basis v_1, \dots, v_n of V . A linear transformation is nonsingular if and only if it sends this basis to another basis of V . Moreover, if w_1, \dots, w_n is any basis of V , by Theorem 6 in Section 10.3 there is a unique linear transformation which sends v_i to w_i , $1 \leq i \leq n$. Thus the number of nonsingular linear transformations from V to itself equals the number of distinct bases of V . This number, which was computed in Example 1 above (with $k = n$), is the order of $GL(V)$.

EXERCISES

1. Let $V = \mathbb{R}^n$ and let (a_1, a_2, \dots, a_n) be a fixed vector in V . Prove that the collection of elements (x_1, x_2, \dots, x_n) of V with $a_1x_1 + a_2x_2 + \dots + a_nx_n = 0$ is a subspace of V . Determine the dimension of this subspace and find a basis.
2. Let V be the collection of polynomials with coefficients in \mathbb{Q} in the variable x of degree at most 5. Prove that V is a vector space over \mathbb{Q} of dimension 6, with $1, x, x^2, \dots, x^5$ as basis. Prove that $1, 1+x, 1+x+x^2, \dots, 1+x+x^2+x^3+x^4+x^5$ is also a basis for V .

3. Let φ be the linear transformation $\varphi : \mathbb{R}^4 \rightarrow \mathbb{R}^1$ such that

$$\begin{aligned}\varphi((1, 0, 0, 0)) &= 1 & \varphi((1, -1, 0, 0)) &= 0 \\ \varphi((1, -1, 1, 0)) &= 1 & \varphi((1, -1, 1, -1)) &= 0.\end{aligned}$$

Determine $\varphi((a, b, c, d))$.

4. Prove that the space of real-valued functions on the closed interval $[a, b]$ is an infinite dimensional vector space over \mathbb{R} , where $a < b$.
5. Prove that the space of continuous real-valued functions on the closed interval $[a, b]$ is an infinite dimensional vector space over \mathbb{R} , where $a < b$.
6. Let V be a vector space of finite dimension. If φ is any linear transformation from V to V prove there is an integer m such that the intersection of the image of φ^m and the kernel of φ^m is $\{0\}$.
7. Let φ be a linear transformation from a vector space V of dimension n to itself that satisfies $\varphi^2 = 0$. Prove that the image of φ is contained in the kernel of φ and hence that the rank of φ is at most $n/2$.
8. Let V be a vector space over F and let φ be a linear transformation of the vector space V to itself. A nonzero element $v \in V$ satisfying $\varphi(v) = \lambda v$ for some $\lambda \in F$ is called an *eigenvector* of φ with *eigenvalue* λ . Prove that for any fixed $\lambda \in F$ the collection of eigenvectors of φ with eigenvalue λ together with 0 forms a subspace of V .
9. Let V be a vector space over F and let φ be a linear transformation of the vector space V to itself. Suppose for $i = 1, 2, \dots, k$ that $v_i \in V$ is an eigenvector for φ with eigenvalue $\lambda_i \in F$ (cf. the preceding exercise) and that all the eigenvalues λ_i are distinct. Prove that v_1, v_2, \dots, v_k are linearly independent. [Use induction on k : write a linear dependence relation among the v_i and apply φ to get another linear dependence relation among the v_i involving the eigenvalues — now subtract a suitable multiple of the first linear relation to get a linear dependence relation on fewer elements.] Conclude that any linear transformation on an n -dimensional vector space has at most n distinct eigenvalues.

In the following exercises let V be a vector space of arbitrary dimension over a field F .

10. Prove that any vector space V has a basis (by convention the null set is the basis for the zero space). [Let \mathcal{S} be the set of subsets of V consisting of linearly independent vectors, partially ordered under inclusion; apply Zorn's Lemma to \mathcal{S} and show a maximal element of \mathcal{S} is a basis.]
11. Refine your argument in the preceding exercise to prove that any set of linearly independent vectors of V is contained in a basis of V .
12. If F is a field with a finite or countable number of elements and V is an infinite dimensional vector space over F with basis \mathcal{B} , prove that the cardinality of V equals the cardinality of \mathcal{B} . Deduce in this case that any two bases of V have the same cardinality.
13. Prove that as vector spaces over \mathbb{Q} , $\mathbb{R}^n \cong \mathbb{R}$, for all $n \in \mathbb{Z}^+$ (note that, in particular, this means \mathbb{R}^n and \mathbb{R} are isomorphic as additive abelian groups).
14. Let \mathcal{A} be a basis for the infinite dimensional space V . Prove that V is isomorphic to the *direct sum* of copies of the field F indexed by the set \mathcal{A} . Prove that the *direct product* of copies of F indexed by \mathcal{A} is a vector space over F and it has strictly larger dimension than the dimension of V (see the exercises in Section 10.3 for the definitions of direct sum and direct product of infinitely many modules).

11.2 THE MATRIX OF A LINEAR TRANSFORMATION

Throughout this section let V, W be vector spaces over the same field F , let $\mathcal{B} = \{v_1, v_2, \dots, v_n\}$ be an (ordered) basis of V , let $\mathcal{E} = \{w_1, w_2, \dots, w_m\}$ be an (ordered) basis of W and let $\varphi \in \text{Hom}(V, W)$ be a linear transformation from V to W . For each $j \in \{1, 2, \dots, n\}$ write the image of v_j under φ in terms of the basis \mathcal{E} :

$$\varphi(v_j) = \sum_{i=1}^m \alpha_{ij} w_i. \quad (11.3)$$

Let $M_{\mathcal{B}}^{\mathcal{E}}(\varphi) = (a_{ij})$ be the $m \times n$ matrix whose i, j entry is α_{ij} (that is, use the coefficients of the w_i 's in the above computation of $\varphi(v_j)$ for the j^{th} column of this matrix). The matrix $M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ is called the *matrix of φ with respect to the bases \mathcal{B}, \mathcal{E}* . The domain basis is the lower and the codomain basis the upper letters appearing after the “*M*.” Given this matrix, we can recover the linear transformation φ as follows: to compute $\varphi(v)$ for $v \in V$, write v in terms of the basis \mathcal{B} :

$$v = \sum_{i=1}^n \alpha_i v_i, \quad \alpha_i \in F,$$

and then calculate the product of the $m \times n$ and $n \times 1$ matrices

$$M_{\mathcal{B}}^{\mathcal{E}}(\varphi) \times \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{pmatrix}.$$

The image of v under φ is given by

$$\varphi(v) = \sum_{i=1}^m \beta_i w_i,$$

i.e., the column vector of coordinates of $\varphi(v)$ with respect to the basis \mathcal{E} are obtained by multiplying the matrix $M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ by the column vector of coordinates of v with respect to the basis \mathcal{B} (sometimes denoted $[\varphi(v)]_{\mathcal{E}} = M_{\mathcal{B}}^{\mathcal{E}}(\varphi)[v]_{\mathcal{B}}$).

Definition. The $m \times n$ matrix $A = (a_{ij})$ associated to the linear transformation φ above is said to *represent* the linear transformation φ with respect to the bases \mathcal{B}, \mathcal{E} . Similarly, φ is the linear transformation represented by A with respect to the bases \mathcal{B}, \mathcal{E} .

Examples

- (1) Let $V = \mathbb{R}^3$ with the standard basis $\mathcal{B} = \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ and let $W = \mathbb{R}^2$ with the standard basis $\mathcal{E} = \{(1, 0), (0, 1)\}$. Let φ be the linear transformation $\varphi(x, y, z) = (x + 2y, x + y + z)$. Since $\varphi(1, 0, 0) = (1, 1)$, $\varphi(0, 1, 0) = (2, 1)$, $\varphi(0, 0, 1) = (0, 1)$, the matrix $A = M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ is the matrix $\begin{pmatrix} 1 & 2 & 0 \\ 1 & 1 & 1 \end{pmatrix}$.

(2) Let $V = W$ be the 2-dimensional space of solutions of the differential equation $y'' - 3y' + 2y = 0$ over \mathbb{C} and let $\mathcal{B} = \mathcal{E}$ be the basis $v_1 = e^t$, $v_2 = e^{2t}$. Since the coefficients of this equation are constants it is easy to check that if y is a solution then its derivative y' is also a solution. It follows that the map $\varphi = d/dt$ = differentiation (with respect to t) is a linear transformation from V to itself. Since $\varphi(v_1) = d(e^t)/dt = e^t = v_1$ and $\varphi(v_2) = d(e^{2t})/dt = 2e^{2t} = 2v_2$ we see that the corresponding matrix with respect to these bases is the diagonal matrix $\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$.

(3) Let $V = W = \mathbb{Q}^3 = \{(x, y, z) \mid x, y, z \in \mathbb{Q}\}$ be the usual 3-dimensional vector space of ordered 3-tuples with entries from the field $F = \mathbb{Q}$ of rational numbers and suppose φ is the linear transformation

$$\varphi(x, y, z) = (9x + 4y + 5z, -4x - 3z, -6x - 4y - 2z), \quad x, y, z \in \mathbb{Q}$$

from V to itself. Take the standard basis $e_1 = (1, 0, 0)$, $e_2 = (0, 1, 0)$, $e_3 = (0, 0, 1)$ for V and for $W = V$. Since $\varphi(1, 0, 0) = (9, -4, -6)$, $\varphi(0, 1, 0) = (4, 0, -4)$, $\varphi(0, 0, 1) = (5, -3, -2)$, the matrix A representing this linear transformation with respect to these bases is

$$A = \begin{pmatrix} 9 & 4 & 5 \\ -4 & 0 & -3 \\ -6 & -4 & -2 \end{pmatrix}.$$

Theorem 10. Let V be a vector space over F of dimension n and let W be a vector space over F of dimension m , with bases \mathcal{B} , \mathcal{E} respectively. Then the map $\text{Hom}_F(V, W) \rightarrow M_{m \times n}(F)$ from the space of linear transformations from V to W to the space of $m \times n$ matrices with coefficients in F defined by $\varphi \mapsto M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ is a vector space isomorphism. In particular, there is a bijective correspondence between linear transformations and their associated matrices with respect to a fixed choice of bases.

Proof: The columns of the matrix $M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ are determined by the action of φ on the basis \mathcal{B} as in equation (3). This shows in particular that the map $\varphi \mapsto M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ is an F -linear map since φ is F -linear. This map is *surjective* since given a matrix M , the map φ defined by equation (3) on a basis and then extended by linearity is a linear transformation with matrix M . The map is *injective* since two linear transformations agreeing on a basis are the same.

Note that different choices of bases give rise to different isomorphisms, so in the same sense that there is no natural choice of basis for a vector space, there is no natural isomorphism between $\text{Hom}_F(V, W)$ and $M_{m \times n}(F)$.

Corollary 11. The dimension of $\text{Hom}_F(V, W)$ is $(\dim V)(\dim W)$.

Proof: The dimension of $M_{m \times n}(F)$ is mn .

Definition. An $m \times n$ matrix A is called *nonsingular* if $Ax = 0$ with $x \in F^n$ implies $x = 0$.

The connection of the term *nonsingular* applied to matrices and to linear transformations is the following: let $A = M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ be the matrix associated to the linear transformation φ (with some choice of bases \mathcal{B} , \mathcal{E}). Then independently of the choice of bases, the $m \times n$ matrix A is nonsingular if and only if the linear transformation φ is a nonsingular linear transformation from the n -dimensional space V to the m -dimensional space W (cf. the exercises).

Assume now that U , V and W are all finite dimensional vector spaces over F with ordered bases \mathcal{D} , \mathcal{B} and \mathcal{E} respectively, where \mathcal{B} and \mathcal{E} are as before and suppose $\mathcal{D} = \{u_1, u_2, \dots, u_k\}$. Assume $\psi : U \rightarrow V$ and $\varphi : V \rightarrow W$ are linear transformations. Their composite, $\varphi \circ \psi$, is a linear transformation from U to W , so we can compute its matrix with respect to the appropriate bases; namely, $M_{\mathcal{D}}^{\mathcal{E}}(\varphi \circ \psi)$ is found by computing

$$\varphi \circ \psi(u_j) = \sum_{i=1}^m \gamma_{ij} w_i$$

and putting the coefficients γ_{ij} down the j^{th} column of $M_{\mathcal{D}}^{\mathcal{E}}(\varphi \circ \psi)$. Next, compute the matrices of ψ and φ separately:

$$\psi(u_j) = \sum_{p=1}^n \alpha_{pj} v_p \quad \text{and} \quad \varphi(v_p) = \sum_{i=1}^m \beta_{ip} w_i$$

so that $M_{\mathcal{D}}^{\mathcal{B}}(\psi) = (\alpha_{pj})$ and $M_{\mathcal{B}}^{\mathcal{E}}(\varphi) = (\beta_{ip})$.

Using these coefficients we can find an expression for the γ 's in terms of the α 's and β 's as follows:

$$\begin{aligned} \varphi \circ \psi(u_j) &= \varphi \left(\sum_{p=1}^n \alpha_{pj} v_p \right) \\ &= \sum_{p=1}^n \alpha_{pj} \varphi(v_p) \\ &= \sum_{p=1}^n \alpha_{pj} \sum_{i=1}^m \beta_{ip} w_i \\ &= \sum_{p=1}^n \sum_{i=1}^m \alpha_{pj} \beta_{ip} w_i. \end{aligned}$$

By interchanging the order of summation in the above double sum we see that γ_{ij} , which is the coefficient of w_i in the above expression, is

$$\gamma_{ij} = \sum_{p=1}^n \alpha_{pj} \beta_{ip}.$$

Computing the product of the matrices for φ and ψ (in that order) we obtain

$$(\beta_{ij})(\alpha_{ij}) = (\delta_{ij}), \quad \text{where} \quad \delta_{ij} = \sum_{p=1}^m \beta_{ip} \alpha_{pj}.$$

By comparing the two sums above and using the commutativity of field multiplication, we see that for all i and j , $\gamma_{ij} = \delta_{ij}$. This computation proves the following result:

Theorem 12. With notations as above, $M_D^E(\varphi \circ \psi) = M_B^E(\varphi)M_D^B(\psi)$, i.e., with respect to a compatible choice of bases, the product of the matrices representing the linear transformations φ and ψ is the matrix representing the composite linear transformation $\varphi \circ \psi$.

Corollary 13. Matrix multiplication is associative and distributive (whenever the dimensions are such as to make products defined). An $n \times n$ matrix A is nonsingular if and only if it is invertible.

Proof: Let A , B and C be matrices such that the products $(AB)C$ and $A(BC)$ are defined, and let S , T and R denote the associated linear transformations. By Theorem 12, the linear transformation corresponding to AB is the composite $S \circ T$ so the linear transformation corresponding to $(AB)C$ is the composite $(S \circ T) \circ R$. Similarly, the linear transformation corresponding to $A(BC)$ is the composite $S \circ (T \circ R)$. Since function composition is associative, these two linear transformations are the same, and so $(AB)C = A(BC)$ by Theorem 10. The distributivity is proved similarly. Note also that it is possible to prove these results by straightforward (albeit tedious) calculations with matrices.

If A is invertible, then $Ax = 0$ implies $x = A^{-1}Ax = A^{-1}0 = 0$, so A is nonsingular. Conversely, if A is nonsingular, fix bases \mathcal{B} , \mathcal{E} for V and let φ be the linear transformation of V to itself represented by A with respect to these bases. By Corollary 9, φ is an isomorphism of V to itself, hence has an inverse, φ^{-1} . Let B be the matrix representing φ^{-1} with respect to the bases \mathcal{E} , \mathcal{B} (note the order). Then $AB = M_B^E(\varphi)M_E^B(\varphi^{-1}) = M_E^E(\varphi \circ \varphi^{-1}) = M_E^E(1) = I$. Similarly, $BA = I$ so B is the inverse of A .

Corollary 14.

- (1) If \mathcal{B} is a basis of the n -dimensional space V , the map $\varphi \mapsto M_B^B(\varphi)$ is a ring and a vector space isomorphism of $\text{Hom}_F(V, V)$ onto the space $M_n(F)$ of $n \times n$ matrices with coefficients in F .
- (2) $GL(V) \cong GL_n(F)$ where $\dim V = n$. In particular, if F is a finite field the order of the finite group $GL_n(F)$ (which equals $|GL(V)|$) is given by the formula at the end of Section 1.

Proof: (1) We have already seen in Theorem 10 that this map is an isomorphism of vector spaces over F . Corollary 13 shows that $M_n(F)$ is a ring under matrix multiplication, and then Theorem 12 shows that multiplication is preserved under this map, hence it is also a ring isomorphism.

(2) This is immediate from (1) since a ring isomorphism sends units to units.

Definition. If A is any $m \times n$ matrix with entries from F , the *row rank* (respectively, *column rank*) of A is the maximal number of linearly independent rows (respectively,

columns) of A (where the rows or columns of A are considered as vectors in affine n -space, m -space, respectively).

The relation between the rank of a matrix and the rank of the associated linear transformation is the following: the rank of φ as a linear transformation equals the column rank of the matrix $M_B^E(\varphi)$ (cf. the exercises). We shall also see that the row rank and the column rank of any matrix are the same.

We now consider the relation of two matrices associated to the same linear transformation of a vector space to itself but with respect to two different choices of bases (cf. the exercises for the general statement regarding a linear transformation from a vector space V to another vector space W).

Definition. Two $n \times n$ matrices A and B are said to be *similar* if there is an invertible (i.e., nonsingular) $n \times n$ matrix P such that $P^{-1}AP = B$. Two linear transformations φ and ψ from a vector space V to itself are said to be *similar* if there is a nonsingular linear transformation ξ from V to V such that $\xi^{-1}\varphi\xi = \psi$.

Suppose \mathcal{B} and \mathcal{E} are two bases of the same vector space V and let $\varphi \in \text{Hom}_F(V, V)$. Let I be the identity map from V to V and let $P = M_{\mathcal{E}}^{\mathcal{B}}(I)$ be its associated matrix (in other words, write the elements of the basis \mathcal{E} in terms of the basis \mathcal{B} — note the order — and use the resulting coordinates for the columns of the matrix P). Note that if $\mathcal{B} \neq \mathcal{E}$ then P is *not* the identity matrix. Then $P^{-1}M_{\mathcal{B}}^{\mathcal{B}}(\varphi)P = M_{\mathcal{E}}^{\mathcal{E}}(\varphi)$. If $[v]_{\mathcal{B}}$ is the $n \times 1$ matrix of coordinates for $v \in V$ with respect to the basis \mathcal{B} , and similarly $[v]_{\mathcal{E}}$ is the $n \times 1$ matrix of coordinates for $v \in V$ with respect to the basis \mathcal{E} , then $[v]_{\mathcal{B}} = P[v]_{\mathcal{E}}$. The matrix P is called the *transition* or *change of basis* matrix from \mathcal{B} to \mathcal{E} and this similarity action on $M_{\mathcal{B}}^{\mathcal{B}}(\varphi)$ is called a *change of basis*. This shows that the matrices associated to the same linear transformation with respect to two different bases are similar.

Conversely, suppose A and B are $n \times n$ matrices similar by a nonsingular matrix P . Let \mathcal{B} be a basis for the n -dimensional vector space V . Define the linear transformation φ of V (with basis \mathcal{B}) to V (again with basis \mathcal{B}) by equation (3) using the given matrix A , i.e.,

$$\varphi(v_j) = \sum_{i=1}^n \alpha_{ij} v_i.$$

Then $A = M_{\mathcal{B}}^{\mathcal{B}}(\varphi)$ by definition of φ . Define a new basis \mathcal{E} of V by using the i^{th} column of P for the coordinates of w_i in terms of the basis \mathcal{B} (so $P = M_{\mathcal{E}}^{\mathcal{B}}(I)$ by definition). Then $B = P^{-1}AP = P^{-1}M_{\mathcal{B}}^{\mathcal{B}}(\varphi)P = M_{\mathcal{E}}^{\mathcal{E}}(\varphi)$ is the matrix associated to φ with respect to the basis \mathcal{E} . This shows that any two similar $n \times n$ matrices arise in this fashion as the matrices representing the same linear transformation with respect to two different choices of bases.

Note that change of basis for a linear transformation from V to itself is the same as conjugation by some element of the group $GL(V)$ of nonsingular linear transformations of V to V . In particular, the relation “similarity” is an equivalence relation whose equivalence classes are the orbits of $GL(V)$ acting by conjugation on $\text{Hom}_F(V, V)$. If

$\varphi \in GL(V)$ (i.e., φ is an invertible linear transformation), then the similarity class of φ is none other than the conjugacy class of φ in the group $GL(V)$.

Example

Let $V = \mathbb{Q}^3$ and let φ be the linear transformation

$$\varphi(x, y, z) = (9x + 4y + 5z, -4x - 3z, -6x - 4y - 2z), \quad x, y, z \in \mathbb{Q}$$

from V to itself we considered in an earlier example. With respect to the standard basis, \mathcal{B} , $b_1 = (1, 0, 0)$, $b_2 = (0, 1, 0)$, $b_3 = (0, 0, 1)$ we saw that the matrix A representing this linear transformation is

$$A = M_{\mathcal{B}}^{\mathcal{B}}(\varphi) = \begin{pmatrix} 9 & 4 & 5 \\ -4 & 0 & -3 \\ -6 & -4 & -2 \end{pmatrix}.$$

Take now the basis, \mathcal{E} , $e_1 = (2, -1, -2)$, $e_2 = (1, 0, -1)$, $e_3 = (3, -2, -2)$ for V (we shall see that this is in fact a basis momentarily). Since

$$\varphi(e_1) = \varphi(2, -1, -2) = (4, -2, -4) = 2 \cdot e_1 + 0 \cdot e_2 + 0 \cdot e_3$$

$$\varphi(e_2) = \varphi(1, 0, -1) = (4, -1, -4) = 1 \cdot e_1 + 2 \cdot e_2 + 0 \cdot e_3$$

$$\varphi(e_3) = \varphi(3, -2, -2) = (9, -6, -6) = 0 \cdot e_1 + 0 \cdot e_2 + 3 \cdot e_3,$$

the matrix representing φ with respect to this basis is the matrix

$$B = M_{\mathcal{E}}^{\mathcal{E}}(\varphi) = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}.$$

Writing the elements of the basis \mathcal{E} in terms of the basis \mathcal{B} we have

$$e_1 = 2b_1 - b_2 - 2b_3$$

$$e_2 = b_1 - b_3$$

$$e_3 = 3b_1 - 2b_2 - 2b_3$$

so the matrix $P = M_{\mathcal{E}}^{\mathcal{B}}(I) = \begin{pmatrix} 2 & 1 & 3 \\ -1 & 0 & -2 \\ -2 & -1 & -2 \end{pmatrix}$ with inverse $P^{-1} = \begin{pmatrix} -2 & -1 & -2 \\ 2 & 2 & 1 \\ 1 & 0 & 1 \end{pmatrix}$

conjugates A into B , i.e., $P^{-1}AP = B$, as can easily be checked. (Note incidentally that since P is invertible this proves that \mathcal{E} is indeed a basis for V .)

We observe in passing that the matrix B representing this linear transformation φ is much simpler than the matrix A representing φ . The study of the simplest possible matrix representing a given linear transformation (and which basis to choose to realize it) is the study of *canonical forms* considered in the next chapter.

Linear Transformations on Tensor Products of Vector Spaces

For convenience we reiterate Corollaries 18 and 19 of Section 10.4 for the special case of vector spaces.

Proposition 15. Let F be a subfield of the field K . If W is an m -dimensional vector space over F with basis w_1, \dots, w_m , then $K \otimes_F W$ is an m -dimensional vector space over K with basis $1 \otimes w_1, \dots, 1 \otimes w_m$.

Proposition 16. Let V and W be finite dimensional vector spaces over the field F with bases v_1, \dots, v_n and w_1, \dots, w_m respectively. Then $V \otimes_F W$ is a vector space over F of dimension nm with basis $v_i \otimes w_j$, $1 \leq i \leq n$ and $1 \leq j \leq m$.

Remark: If v and w are nonzero elements of V and W , respectively, then it follows from the proposition that $v \otimes w$ is a nonzero element of $V \otimes_F W$, because we may always build bases of V and W whose first basis vectors are v , w , respectively. In a tensor product $M \otimes_R N$ of two R -modules where R is not a field it is in general substantially more difficult to determine when the tensor product $m \otimes n$ of two nonzero elements is zero.

Now let V, W, X, Y be finite dimensional vector spaces over F and let

$$\varphi : V \rightarrow X \quad \text{and} \quad \psi : W \rightarrow Y$$

be linear transformations. We compute a matrix of the linear transformation

$$\varphi \otimes \psi : V \otimes W \rightarrow X \otimes Y.$$

Let $\mathcal{B}_1 = \{v_1, \dots, v_n\}$ and $\mathcal{B}_2 = \{w_1, \dots, w_m\}$ be (ordered) bases of V and W respectively, and let $\mathcal{E}_1 = \{x_1, \dots, x_r\}$ and $\mathcal{E}_2 = \{y_1, \dots, y_s\}$ be (ordered) bases of X and Y respectively. Let $\mathcal{B} = \{v_i \otimes w_j\}$ and $\mathcal{E} = \{x_i \otimes y_j\}$ be the bases of $V \otimes W$ and $X \otimes Y$ given by Proposition 16; we shall order these shortly. Suppose

$$\varphi(v_i) = \sum_{p=1}^r \alpha_{pi} x_p \quad \text{and} \quad \psi(w_j) = \sum_{q=1}^s \beta_{qj} y_q.$$

Then

$$\begin{aligned} (\varphi \otimes \psi)(v_i \otimes w_j) &= (\varphi(v_i)) \otimes (\psi(w_j)) \\ &= \left(\sum_{p=1}^r \alpha_{pi} x_p \right) \otimes \left(\sum_{q=1}^s \beta_{qj} y_q \right) \\ &= \sum_{p=1}^r \sum_{q=1}^s \alpha_{pi} \beta_{qj} (x_p \otimes y_q). \end{aligned} \tag{11.8}$$

In view of the order of summation in (11.8) we order the basis \mathcal{E} into r ordered sets, with the p^{th} list being $x_p \otimes y_1, x_p \otimes y_2, \dots, x_p \otimes y_s$, and similarly order the basis \mathcal{B} . Then equation (8) determines the column entries for the corresponding matrix of $\varphi \otimes \psi$. The resulting matrix $M_{\mathcal{B}}^{\mathcal{E}}(\varphi \otimes \psi)$ is an $r \times n$ block matrix whose p, q block is the $s \times m$ matrix $\alpha_{p,q} M_{\mathcal{B}_2}^{\mathcal{E}_2}(\psi)$. In other words, the matrix for $\varphi \otimes \psi$ is obtained by taking the matrix for φ and multiplying each entry by the matrix for ψ . Such matrices have a name:

Definition. Let $A = (\alpha_{ij})$ and B be $r \times n$ and $s \times m$ matrices, respectively, with coefficients from any commutative ring. The *Kronecker product* or *tensor product* of A and B , denoted by $A \otimes B$, is the $rs \times nm$ matrix consisting of an $r \times n$ block matrix whose i, j block is the $s \times m$ matrix $\alpha_{ij} B$.

With this terminology we have

Proposition 17. Let $\varphi : V \rightarrow X$ and $\psi : W \rightarrow Y$ be linear transformations of finite dimensional vector spaces. Then the Kronecker product of matrices representing φ and ψ is a matrix representation of $\varphi \otimes \psi$.

Example

Let $V = X = \mathbb{R}^3$, both with basis v_1, v_2, v_3 , and $W = Y = \mathbb{R}^2$, both with basis w_1, w_2 . Suppose $\varphi : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is the linear transformation given by $\varphi(av_1 + bv_2 + cv_3) = cv_1 + 2av_2 - 3bv_3$ and $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is the linear transformation given by $\psi(aw_1 + bw_2) = (a+3b)w_1 + (4b-2a)w_2$. With respect to the chosen bases, the matrices for φ and ψ are

$$\begin{pmatrix} 0 & 0 & 1 \\ 2 & 0 & 0 \\ 0 & -3 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 & 3 \\ -2 & 4 \end{pmatrix},$$

respectively. Then with respect to the ordered basis

$$\mathcal{B} = \{v_1 \otimes w_1, v_1 \otimes w_2, v_2 \otimes w_1, v_2 \otimes w_2, v_3 \otimes w_1, v_3 \otimes w_2\}$$

we have

$$M_{\mathcal{B}}^{\mathcal{B}}(\varphi \otimes \psi) = \left(\begin{array}{cc|cc|cc} 0 & 0 & 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 0 & -2 & 4 \\ \hline 2 & 6 & 0 & 0 & 0 & 0 \\ -4 & 8 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & -3 & -9 & 0 & 0 \\ 0 & 0 & 6 & -12 & 0 & 0 \end{array} \right),$$

obtained (as indicated by the dashed lines) by multiplying the 2×2 matrix for ψ successively by the entries in the matrix for φ .

EXERCISES

- Let V be the collection of polynomials with coefficients in \mathbb{Q} in the variable x of degree at most 5. Determine the transition matrix from the basis $1, x, x^2, \dots, x^5$ for V to the basis $1, 1+x, 1+x+x^2, \dots, 1+x+x^2+x^3+x^4+x^5$ for V .
- Let V be the vector space of the preceding exercise. Let $\varphi = d/dx$ be the linear transformation of V to itself given by usual differentiation of a polynomial with respect to x . Determine the matrix of φ with respect to the two bases for V in the previous exercise.
- Let V be the collection of polynomials with coefficients in F in the variable x of degree at most n . Determine the transition matrix from the basis $1, x, x^2, \dots, x^n$ for V to the elements

$$1, x - \lambda, \dots, (x - \lambda)^{n-1}, (x - \lambda)^n$$

where λ is a fixed element of F . Conclude that these elements are a basis for V .

- Let φ be the linear transformation of \mathbb{R}^2 to itself given by rotation counterclockwise around the origin through an angle θ . Show that the matrix of φ with respect to the standard basis for \mathbb{R}^2 is $\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$.
- Show that the $m \times n$ matrix A is nonsingular if and only if the linear transformation φ is a nonsingular linear transformation from the n -dimensional space V to the m -dimensional space W , where $A = M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$, regardless of the choice of bases \mathcal{B} and \mathcal{E} .

6. Prove if $\varphi \in \text{Hom}_F(F^n, F^m)$, and \mathcal{B}, \mathcal{E} are the natural bases of F^n, F^m respectively, then the range of φ equals the span of the set of columns of $M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$. Deduce that the rank of φ (as a linear transformation) equals the column rank of $M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$.
7. Prove that any two similar matrices have the same row rank and the same column rank.
8. Let V be an n -dimensional vector space over F and let φ be a linear transformation of the vector space V to itself.
- Prove that if V has a basis consisting of eigenvectors for φ (cf. Exercise 8 of Section 1) then the matrix representing φ with respect to this basis (for both domain and range) is diagonal with the eigenvalues as diagonal entries.
 - If A is the $n \times n$ matrix representing φ with respect to a given basis for V (for both domain and range) prove that A is similar to a diagonal matrix if and only if V has a basis of eigenvectors for φ .
9. If W is a subspace of the vector space V stable under the linear transformation φ (i.e., $\varphi(W) \subseteq W$), show that φ induces linear transformations $\varphi|_W$ on W and $\tilde{\varphi}$ on the quotient vector space V/W . If $\varphi|_W$ and $\tilde{\varphi}$ are nonsingular prove φ is nonsingular. Prove the converse holds if V has finite dimension and give a counterexample with V infinite dimensional.
10. Let V be an n -dimensional vector space and let φ be a linear transformation of V to itself. Suppose W is a subspace of V of dimension m that is stable under φ .
- Prove that there is a basis for V with respect to which the matrix for φ is of the form
- $$\begin{pmatrix} A & B \\ 0 & C \end{pmatrix}$$
- where A is an $m \times m$ matrix, B is an $m \times (n-m)$ matrix and C is an $(n-m) \times (n-m)$ matrix (such a matrix is called *block upper triangular*).
- Prove that if there is a subspace W' invariant under φ so that $V = W \oplus W'$ decomposes as a direct sum then the bases for W and W' give a basis for V with respect to which the matrix for φ is *block diagonal*:
- $$\begin{pmatrix} A & 0 \\ 0 & C \end{pmatrix}$$
- where A is an $m \times m$ matrix and C is an $(n-m) \times (n-m)$ matrix.
- Prove conversely that if there is a basis for V with respect to which φ is block diagonal as in (b) then there are φ -invariant subspaces W and W' of dimensions m and $n-m$, respectively, with $V = W \oplus W'$.
11. Let φ be a linear transformation from the finite dimensional vector space V to itself such that $\varphi^2 = \varphi$.
- Prove that $\text{image } \varphi \cap \ker \varphi = 0$.
 - Prove that $V = \text{image } \varphi \oplus \ker \varphi$.
 - Prove that there is a basis of V such that the matrix of φ with respect to this basis is a diagonal matrix whose entries are all 0 or 1.
- A linear transformation φ satisfying $\varphi^2 = \varphi$ is called an *idempotent* linear transformation. This exercise proves that idempotent linear transformations are simply projections onto some subspace.
12. Let $V = \mathbb{R}^2$, $v_1 = (1, 0)$, $v_2 = (0, 1)$, so that v_1, v_2 are a basis for V . Let φ be the linear transformation of V to itself whose matrix with respect to this basis is $\begin{pmatrix} 2 & 1 \\ 0 & 2 \end{pmatrix}$. Prove that if W is the subspace generated by v_1 then W is stable under the action of φ . Prove that there is no subspace W' invariant under φ so that $V = W \oplus W'$.

13. Let V be a vector space of dimension n and let W be a vector space of dimension m over a field F . Suppose A is the $m \times n$ matrix representing a linear transformation φ from V to W with respect to the bases \mathcal{B}_1 for V and \mathcal{E}_1 for W . Suppose similarly that B is the $m \times n$ matrix representing φ with respect to the bases \mathcal{B}_2 for V and \mathcal{E}_2 for W . Let $P = M_{\mathcal{B}_2}^{\mathcal{B}_1}(I)$ where I denotes the identity map from V to V , and let $Q = M_{\mathcal{E}_2}^{\mathcal{E}_1}(I)$ where I denotes the identity map from W to W . Prove that $Q^{-1} = M_{\mathcal{E}_1}^{\mathcal{E}_2}(I)$ and that $Q^{-1}AP = B$, giving the general relation between matrices representing the same linear transformation but with respect to different choices of bases.

The following exercises recall the *Gauss–Jordan* elimination process. This is one of the fastest computational methods for the solution of a number of problems involving vector spaces — solving systems of linear equations, determining inverses of matrices, computing determinants, determining the span of a set of vectors, determining linear independence of a set of vectors etc.

Consider the system of m linear equations

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= c_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= c_2 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n &= c_m \end{aligned} \tag{11.4}$$

in the n unknowns x_1, x_2, \dots, x_n where $a_{ij}, c_i, i = 1, 2, \dots, m, j = 1, 2, \dots, n$ are elements of the field F . Associated to this system is the *coefficient matrix*:

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}$$

and the *augmented matrix*:

$$(A | C) = \left(\begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & c_1 \\ a_{21} & a_{22} & \dots & a_{2n} & c_2 \\ \vdots & \vdots & & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} & c_m \end{array} \right)$$

(the term *augmented* refers to the presence of the column matrix $C = (c_i)$ in addition to the coefficient matrix $A = (a_{ij})$). The set of solutions in F of this system of equations is not altered if we perform any of the following three operations:

- (1) interchange any two equations
- (2) add a multiple of one equation to another
- (3) multiply any equation by a nonzero element from F ,

which correspond to the following three *elementary row operations* on the augmented matrix:

- (1) interchange any two rows
- (2) add a multiple of one row to another
- (3) multiply any row by a unit in F , i.e., by any nonzero element in F .

If a matrix A can be transformed into a matrix C by a series of elementary row operations then A is said to be *row reduced* to C .

- 14.** Prove that if A can be row reduced to C then C can be row reduced to A . Prove that the relation “ $A \sim C$ if and only if A can be row reduced to C ” is an equivalence relation. [Observe that the elementary row operations are reversible.]

Matrices lying in the same equivalence class under this equivalence relation are said to be *row equivalent*.

- 15.** Prove that the row rank of two row equivalent matrices is the same. [It suffices to prove this for two matrices differing by an elementary row operation.]

An $m \times n$ matrix is said to be in *reduced row echelon form* if

- (a) the first nonzero entry a_{ij_i} in row i is 1 and all other entries in the corresponding j_i^{th} column are zero, and
- (b) $j_1 < j_2 < \dots < j_r$, where r is the number of nonzero rows, i.e., the number of initial zeros in each row is strictly increasing (hence the term *echelon*).

An augmented matrix $(A | C)$ is said to be in reduced row echelon form if its coefficient matrix A is in reduced row echelon form. For example, the following two matrices are in reduced row echelon form:

$$\left(\begin{array}{ccccc|c} 1 & 0 & 5 & 7 & 0 & 3 \\ 0 & 1 & -1 & 1 & 0 & -4 \\ 0 & 0 & 0 & 0 & 1 & 6 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right) \quad \left(\begin{array}{ccccc|c} 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 & -3 \end{array} \right)$$

(with $j_1 = 1$, $j_2 = 2$, $j_3 = 5$ for the first matrix and $j_1 = 2$, $j_2 = 4$ for the second matrix). The first nonzero entry in any given row of the coefficient matrix of a reduced row echelon augmented matrix (in position (i, j_i) by definition) is sometimes referred to as a *pivotal element* (so the pivotal elements in the first matrix are in positions $(1,1)$, $(2,2)$ and $(3,5)$ and the pivotal elements in the second matrix are in positions $(1,2)$ and $(2,4)$). The columns containing pivotal elements will be called *pivotal columns* and the columns of the coefficient matrix not containing pivotal elements will be called *nonpivotal*.

- 16.** Prove by induction that any augmented matrix can be put in reduced row echelon form by a series of elementary row operations.
17. Let A and C be two matrices in reduced row echelon form. Prove that if A and C are row equivalent then $A = C$.
18. Prove that the row rank of a matrix in reduced row echelon form is the number of nonzero rows.
19. Prove that the reduced row echelon forms of the matrices

$$\left(\begin{array}{ccccc|c} 1 & 1 & 4 & 8 & 0 & -1 \\ 1 & 2 & 3 & 9 & 0 & -5 \\ 0 & -2 & 2 & -2 & 1 & 14 \\ 1 & 4 & 1 & 11 & 0 & -13 \end{array} \right) \quad \left(\begin{array}{ccccc|c} 0 & -3 & 3 & 1 & 5 \\ 0 & 1 & -1 & 0 & 0 \\ 0 & 2 & -2 & 0 & -3 \end{array} \right)$$

are the two matrices preceding Exercise 16.

The point of the reduced row echelon form is that the corresponding system of linear equations is in a particularly simple form, from which the solutions to the system $AX = C$ in (4) can be determined immediately:

- 20.** (*Solving Systems of Linear Equations*) Let $(A' | C')$ be the reduced row echelon form of the augmented matrix $(A | C)$. The number of zero rows of A' is clearly at least as great as the number of zero rows of $(A' | C')$.

- (a) Prove that if the number of zero rows of A' is strictly larger than the number of zero rows of $(A' \mid C')$ then there are no solutions to $AX = C$.

By (a) we may assume that A' and $(A' \mid C')$ have the same number, r , of nonzero rows (so $n \geq r$).

- (b) Prove that if $r = n$ then there is precisely one solution to the system of equations $AX = C$.
(c) Prove that if $r < n$ then there are infinitely many solutions to the system of equations $AX = C$. Prove in fact that the values of the $n - r$ variables corresponding to the nonpivot columns of $(A' \mid C')$ can be chosen arbitrarily and that the remaining r variables corresponding to the pivotal columns of $(A' \mid C')$ are then determined uniquely.

21. Determine the solutions of the following systems of equations:

(a)

$$\begin{array}{rcl} -3x + 3y + z & = & 5 \\ x - y & = & 0 \\ 2x - 2y & = & -3 \end{array}$$

(b)

$$\begin{array}{rcl} x - 2y + z & = & 5 \\ x - 4y - 6z & = & 10 \\ 4x - 11y + 11z & = & 12 \end{array}$$

(c)

$$\begin{array}{rcl} x - 2y + z & = & 5 \\ y - 2z & = & 17 \\ 2x - 3y & = & 27 \end{array}$$

(d)

$$\begin{array}{rcl} x + y - 3z + 2u & = & 2 \\ 3x - 2y + 5z + u & = & 1 \\ 6x + y - 4z + 3u & = & 7 \\ 2x + 2y - 6z & = & 4 \end{array}$$

(e)

$$\begin{array}{rcl} x + y + 4z + 8u & - w & = -1 \\ x + 2y + 3z + 9u & - 5w & = -2 \\ -2y + 2z - 2u + v + 14w & = & 3 \\ x + 4y + z + 11u & - 13w & = -4 \end{array}$$

22. Suppose A and B are two row equivalent $m \times n$ matrices.

- (a) Prove that the set

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

of solutions to the homogeneous linear equations $AX = 0$ as in equation (4) above are the same as the set of solutions to the homogeneous linear equations $BX = 0$. [It suffices to prove this for two matrices differing by an elementary row operation.]

- (b) Prove that any linear dependence relation satisfied by the columns of A viewed as vectors in F^m is also satisfied by the columns of B .

- (c) Conclude from (b) that the number of linearly independent columns of A is the same as the number of linearly independent columns of B .

23. Let A' be a matrix in reduced row echelon form.

- (a) Prove that the nonzero rows of A' are linearly independent. Prove that the pivotal columns of A' are linearly independent and that the nonpivotal columns of A' are linearly dependent on the pivotal columns. (Note the role the pivotal elements play.)
 (b) Prove that the number of linearly independent columns of a matrix in reduced row echelon form is the same as the number of linearly independent rows, i.e., the row rank and the column rank of such a matrix are the same.

24. Use the previous two exercises and Exercise 15 above to prove in general that the row rank and the column rank of a matrix are the same.

25. (*Computing Inverses of Matrices*) Let A be an $n \times n$ matrix.

- (a) Show that A has an inverse matrix B with columns B_1, B_2, \dots, B_n if and only if the systems of equations:

$$AB_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}, \quad AB_2 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{pmatrix}, \quad \dots, \quad AB_n = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}$$

have solutions.

- (b) Prove that A has an inverse if and only if A is row equivalent to the $n \times n$ identity matrix.
 (c) Prove that A has an inverse B if and only if the augmented matrix $(A | I)$ can be row reduced to the augmented matrix $(I | B)$ where I is the $n \times n$ identity matrix.

26. Determine the inverses of the following matrices using row reduction:

$$A = \begin{pmatrix} -7 & -1 & -4 \\ 7 & 1 & 3 \\ 1 & 0 & 0 \end{pmatrix} \quad B = \begin{pmatrix} 1 & 1 & 0 & 2 \\ 0 & 2 & 1 & -1 \\ 0 & 2 & 0 & 0 \\ -1 & 1 & 1 & 0 \end{pmatrix}.$$

27. (*Computing Spans, Linear Independence and Linear Dependencies in Vector Spaces*) Let V be an m -dimensional vector space with basis e_1, e_2, \dots, e_m and let v_1, v_2, \dots, v_n be vectors in V . Let A be the $m \times n$ matrix whose columns are the coordinates of the vectors v_i (with respect to the basis e_1, e_2, \dots, e_m) and let A' be the reduced row echelon form of A .

- (a) Let B be any matrix row equivalent to A . Let w_1, w_2, \dots, w_n be the vectors whose coordinates (with respect to the basis e_1, e_2, \dots, e_m) are the columns of B . Prove that any linear relation

$$x_1 v_1 + x_2 v_2 + \dots + x_n v_n = 0 \tag{11.5}$$

satisfied by v_1, v_2, \dots, v_n is also satisfied when v_i is replaced by w_i , $i = 1, 2, \dots, n$.

- (b) Prove that the vectors whose coordinates are given by the pivotal columns of A' are linearly independent and that the vectors whose coordinates are given by the nonpivotal columns of A' are linearly dependent on these.
 (c) (*Determining Linear Independence of Vectors*) Prove that the vectors v_1, v_2, \dots, v_n are linearly independent if and only if A' has n nonzero rows (i.e., has rank n).
 (d) (*Determining Linear Dependencies of Vectors*) By (c), the vectors v_1, v_2, \dots, v_n are linearly dependent if and only if A' has nonpivotal columns. The solutions to (5)

defining linear dependence relations among v_1, v_2, \dots, v_n are given by the linear equations defined by A' . Show that each of the variables x_1, x_2, \dots, x_n in (5) corresponding to the nonpivot columns of A' can be prescribed arbitrarily and the values of the remaining variables are then uniquely determined to give a linear dependence relation among v_1, v_2, \dots, v_n as in (5).

- (e) (*Determining the Span of a Set of Vectors*) Prove that the subspace W spanned by v_1, v_2, \dots, v_n has dimension r where r is the number of nonzero rows of A' and that a basis for W is given by the original vectors v_{j_i} ($i = 1, 2, \dots, r$) corresponding to the pivotal columns of A' .

28. Let $V = \mathbb{R}^5$ with the standard basis and consider the vectors

$$v_1 = (1, 1, 3, -2, 3), \quad v_2 = (0, 1, 0, -1, 0), \quad v_3 = (2, 3, 6, -5, 6)$$

$$v_4 = (0, 3, 1, -3, 1), \quad v_5 = (2, -1, -1, -1, -1).$$

- (a) Show that the reduced row echelon form of the matrix

$$A = \begin{pmatrix} 1 & 0 & 2 & 0 & 2 \\ 1 & 1 & 3 & 3 & -1 \\ 3 & 0 & 6 & 1 & -1 \\ -2 & -1 & -5 & -3 & -1 \\ 3 & 0 & 6 & 1 & -1 \end{pmatrix}$$

whose columns are the coordinates of v_1, v_2, v_3, v_4, v_5 is the matrix

$$A' = \begin{pmatrix} 1 & 0 & 2 & 0 & 2 \\ 0 & 1 & 1 & 0 & 18 \\ 0 & 0 & 0 & 1 & -7 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

where the 1st, 2nd and 4th columns are pivotal and the remaining two are nonpivotal.

- (b) Conclude that these vectors are linearly dependent, that the subspace W spanned by v_1, v_2, v_3, v_4, v_5 is 3-dimensional and that the vectors

$$v_1 = (1, 1, 3, -2, 3), \quad v_2 = (0, 1, 0, -1, 0) \quad \text{and} \quad v_4 = (0, 3, 1, -3, 1)$$

are a basis for W .

- (c) Conclude from (a) that the coefficients x_1, x_2, x_3, x_4, x_5 of any linear relation

$$x_1 v_1 + x_2 v_2 + x_3 v_3 + x_4 v_4 + x_5 v_5 = 0$$

satisfied by v_1, v_2, v_3, v_4, v_5 are given by the equations

$$\begin{aligned} x_1 &+ 2x_3 &+ 2x_5 &= 0 \\ x_2 &+ x_3 &+ 18x_5 &= 0 \\ x_4 &- 7x_5 &= 0. \end{aligned}$$

Deduce that the 3rd and 5th variables, namely x_3 and x_5 , corresponding to the non-pivotal columns of A' , can be prescribed arbitrarily and the remaining variables are then uniquely determined as:

$$x_1 = -2x_3 - 2x_5$$

$$x_2 = -x_3 - 18x_5$$

$$x_4 = 7x_5$$

to give all the linear dependence relations satisfied by v_1, v_2, v_3, v_4, v_5 . In particular show that

$$-2v_1 - v_2 + v_3 = 0$$

and

$$-2v_1 - 18v_2 + 7v_4 + v_5 = 0$$

corresponding to $(x_3 = 1, x_5 = 0)$ and $(x_3 = 0, x_5 = 1)$, respectively.

- 29.** For each exercise below, determine whether the given vectors in \mathbb{R}^4 are linearly independent. If they are linearly dependent, determine an explicit linear dependence among them.
- (a) $(1, -4, 3, 0), (0, -1, 4, -3), (1, -1, 1, -1), (2, 2, -1, -3)$.
 - (b) $(1, -2, 4, 1), (2, -3, 9, -1), (1, 0, 6, -5), (2, -5, 7, 5)$.
 - (c) $(1, -2, 0, 1), (2, -2, 0, 0), (-1, 3, 0, -2), (-2, 1, 0, 1)$.
 - (d) $(0, 1, 1, 0), (1, 0, 1, 1), (2, 2, 2, 0), (0, -1, 1, 1)$.
- 30.** For each exercise below, determine the subspace spanned in \mathbb{R}^4 by the given vectors and give a basis for this subspace.
- (a) $(1, -2, 5, 3), (2, 3, 1, -4), (3, 8, -3, -5)$.
 - (b) $(2, -5, 3, 0), (0, -2, 5, -3), (1, -1, 1, -1), (-3, 2, -1, 2)$.
 - (c) $(1, -2, 0, 1), (2, -2, 0, 0), (-1, 3, 0, -2), (-2, 1, 0, 1)$.
 - (d) $(1, 1, 0, -1), (1, 2, 3, 0), (2, 3, 3, -1), (1, 2, 2, -2), (2, 3, 2, -3), (1, 3, 4, -3)$.

- 31.** (*Computing the Image and Kernel of a Linear Transformation*) Let V be an n -dimensional vector space with basis e_1, e_2, \dots, e_n and let W be an m -dimensional vector space with basis f_1, f_2, \dots, f_m . Let φ be a linear transformation from V to W and let A be the corresponding $m \times n$ matrix with respect to these bases: $A = (a_{ij})$ where

$$\varphi(e_j) = \sum_{i=1}^m a_{ij} f_i, \quad j = 1, 2, \dots, n,$$

i.e., the columns of A are the coordinates of the vectors $\varphi(e_1), \varphi(e_2), \dots, \varphi(e_n)$ with respect to the basis f_1, f_2, \dots, f_m of W . Let A' be the reduced row echelon form of A .

- (a) (*Determining the Image of a Linear Transformation*) Prove that the image $\varphi(V)$ of V under φ has dimension r where r is the number of nonzero rows of A' and that a basis for $\varphi(V)$ is given by the vectors $\varphi(e_{j_i})$ ($i = 1, 2, \dots, r$), i.e., the columns of A corresponding to the pivotal columns of A' give the coordinates of a basis for the image of φ .
- (b) (*Determining the Kernel of a Linear Transformation*) The elements in the kernel of φ are the vectors in V whose coordinates (x_1, x_2, \dots, x_n) with respect to the basis e_1, e_2, \dots, e_n satisfy the equation

$$A \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = 0,$$

and the solutions x_1, x_2, \dots, x_n to this system of linear equations are determined by the matrix A' .

- (i) Prove that φ is injective if and only if A' has n nonzero rows (i.e., has rank n).
- (ii) By (i), the kernel of φ is nontrivial if and only if A' has nonpivotal columns. Show that each of the variables x_1, x_2, \dots, x_n above corresponding to the nonpivotal columns of A' can be prescribed arbitrarily and the values of the remaining variables are then

uniquely determined to give an element $x_1e_1 + x_2e_2 + \dots + x_ne_n$ in the kernel of φ . In particular, show that the coordinates of a basis for the kernel are obtained by successively setting one nonpivotal variable equal to 1 and all other nonpivotal variables to 0 and solving for the remaining pivotal variables. Conclude that the kernel of φ has dimension $n - r$ where r is the rank of A .

32. Let $V = \mathbb{R}^5$ and $W = \mathbb{R}^4$ with the standard bases. Let φ be the linear transformation $\varphi : V \rightarrow W$ defined by

$$\varphi(x, y, z, u, v) = (x + 2y + 3z + 4u + 4v, -2x - 4y + 2v, x + 2y + u - 2v, x + 2y - v).$$

- (a) Prove that the matrix A corresponding to φ and these bases is

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 & 4 \\ -2 & -4 & 0 & 0 & 2 \\ 1 & 2 & 0 & 1 & -2 \\ 1 & 2 & 0 & 0 & -1 \end{pmatrix}$$

and that the reduced row echelon matrix A' row equivalent to A is

$$A' = \begin{pmatrix} 1 & 2 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 & 3 \\ 0 & 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

where the 1st, 3rd and 4th columns are pivotal and the remaining two are nonpivotal.

- (b) Conclude that the image of φ is 3-dimensional and that the image of the 1st, 3rd and 4th basis elements of V , namely, $(1, -2, 1, 1)$, $(3, 0, 0, 0)$ and $(4, 0, 1, 0)$ give a basis for the image $\varphi(V)$ of V .
(c) Conclude from (a) that the elements in the kernel of φ are the vectors (x, y, z, u, v) satisfying the equations

$$\begin{aligned} x + 2y - v &= 0 \\ z + 3v &= 0 \\ u - v &= 0. \end{aligned}$$

Deduce that the 2nd and 5th variables, namely y and v , corresponding to the nonpivotal columns of A' can be prescribed arbitrarily and the remaining variables are then uniquely determined as

$$\begin{aligned} x &= -2y + v \\ z &= -3v \\ u &= v. \end{aligned}$$

Show that $(-2, 1, 0, 0, 0)$ and $(1, 0, -3, 1, 1)$ give a basis for the 2-dimensional kernel of φ , corresponding to $(y = 1, v = 0)$ and $(y = 0, v = 1)$, respectively.

33. Let φ be the linear transformation from \mathbb{R}^4 to itself defined by the matrix

$$A = \begin{pmatrix} 1 & -1 & 0 & 3 \\ -1 & 2 & 1 & -1 \\ -1 & 1 & 0 & -3 \\ 1 & -2 & -1 & 1 \end{pmatrix}$$

with respect to the standard basis for \mathbb{R}^4 . Determine a basis for the image and for the kernel of φ .

34. Let φ be the linear transformation $\varphi : \mathbb{R}^4 \rightarrow \mathbb{R}^2$ such that

$$\begin{aligned}\varphi((1, 0, 0, 0)) &= (1, -1) & \varphi((1, -1, 0, 0)) &= (0, 0) \\ \varphi((1, -1, 1, 0)) &= (1, -1) & \varphi((1, -1, 1, -1)) &= (0, 0).\end{aligned}$$

Determine a basis for the image and for the kernel of φ .

35. Let V be the set of all 2×2 matrices with real entries and let $\varphi : V \rightarrow \mathbb{R}$ be the map defined by sending a matrix $A \in V$ to the sum of the diagonal entries of A (the *trace* of A).

- (a) Show that

$$\left(\begin{array}{cc} 1 & 0 \\ 0 & 0 \end{array}\right), \quad \left(\begin{array}{cc} 0 & 1 \\ 0 & 0 \end{array}\right), \quad \left(\begin{array}{cc} 0 & 0 \\ 1 & 0 \end{array}\right), \quad \left(\begin{array}{cc} 0 & 0 \\ 0 & 1 \end{array}\right)$$

is a basis for V .

- (b) Prove that φ is a linear transformation and determine the matrix of φ with respect to the basis in (a) for V . Determine the dimension of and a basis for the kernel of φ .

36. Let V be the 6-dimensional vector space over \mathbb{Q} consisting of the polynomials in the variable x of degree at most 5. Let φ be the map of V to itself defined by $\varphi(f) = x^2 f'' - 6x f' + 12f$, where f'' denotes the usual second derivative (with respect to x) of the polynomial $f \in V$ and f' similarly denotes the usual first derivative.

- (a) Prove that φ is a linear transformation of V to itself.

- (b) Determine a basis for the image and for the kernel of φ .

37. Let V be the 7-dimensional vector space over the field F consisting of the polynomials in the variable x of degree at most 6. Let φ be the linear transformation of V to itself defined by $\varphi(f) = f'$, where f' denotes the usual derivative (with respect to x) of the polynomial $f \in V$. For each of the fields below, determine a basis for the image and for the kernel of φ :

- (a) $F = \mathbb{R}$

- (b) $F = \mathbb{F}_2$, the finite field of 2 elements (note that, for example, $(x^2)' = 2x = 0$ over this field)

- (c) $F = \mathbb{F}_3$

- (d) $F = \mathbb{F}_5$.

38. Let A and B be square matrices. Prove that the trace of their Kronecker product is the product of their traces: $\text{tr}(A \otimes B) = \text{tr}(A) \text{tr}(B)$. (Recall that the trace of a square matrix is the sum of its diagonal entries.)

39. Let F be a subfield of K and let $\psi : V \rightarrow W$ be a linear transformation of finite dimensional vector spaces over F .

- (a) Prove that $1 \otimes \psi$ is a K -linear transformation from the vector spaces $K \otimes_F V$ to $K \otimes_F W$ over K . (Here 1 denotes the identity map from K to itself.)

- (b) Let $\mathcal{B} = \{v_1, \dots, v_n\}$ and $\mathcal{E} = \{w_1, \dots, w_m\}$ be bases of V and W respectively. Prove that the matrix of $1 \otimes \psi$ with respect to the bases $\{1 \otimes v_1, \dots, 1 \otimes v_n\}$ and $\{1 \otimes w_1, \dots, 1 \otimes w_m\}$ is the same as the matrix of ψ with respect to \mathcal{B} and \mathcal{E} .

11.3 DUAL VECTOR SPACES

Definition.

- (1) For V any vector space over F let $V^* = \text{Hom}_F(V, F)$ be the space of linear transformations from V to F , called the *dual space* of V . Elements of V^* are called *linear functionals*.

- (2) If $\mathcal{B} = \{v_1, v_2, \dots, v_n\}$ is a basis of the finite dimensional space V , define $v_i^* \in V^*$ for each $i \in \{1, 2, \dots, n\}$ by its action on the basis \mathcal{B} :

$$v_i^*(v_j) = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad 1 \leq j \leq n. \quad (11.6)$$

Proposition 18. With notations as above, $\{v_1^*, v_2^*, \dots, v_n^*\}$ is a basis of V^* . In particular, if V is finite dimensional then V^* has the same dimension as V .

Proof. Observe that since V is finite dimensional, $\dim V^* = \dim \text{Hom}_F(V, F) = \dim V = n$ (Corollary 11), so since there are n of the v_i^* 's it suffices to prove that they are linearly independent. If

$$\alpha_1 v_1^* + \alpha_2 v_2^* + \cdots + \alpha_n v_n^* = 0 \quad \text{in } \text{Hom}_F(V, F),$$

then applying this element to v_i and using equation (6) above we obtain $\alpha_i = 0$. Since i is arbitrary these elements are linearly independent.

Definition. The basis $\{v_1^*, v_2^*, \dots, v_n^*\}$ of V^* is called the *dual basis* to $\{v_1, v_2, \dots, v_n\}$.

The exercises later show that if V is infinite dimensional it is always true that $\dim V < \dim V^*$. For spaces of arbitrary dimension the space V^* is the “algebraic” dual space to V . If V has some additional structure, for example a continuous structure (i.e., a topology), then one may define other types of dual spaces (e.g., the continuous dual of V , defined by requiring the linear functionals to be *continuous* maps). One has to be careful when reading other works (particularly analysis books) to ascertain what qualifiers are implicit in the use of the terms “dual space” and “linear functional.”

Example

Let $[a, b]$ be a closed interval in \mathbb{R} and let V be the real vector space of all continuous functions $f : [a, b] \rightarrow \mathbb{R}$. If $a < b$, V is infinite dimensional. For each $g \in V$ the function $\varphi_g : V \rightarrow \mathbb{R}$ defined by $\varphi_g(f) = \int_a^b f(t)g(t)dt$ is a linear functional on V .

Definition. The dual of V^* , namely V^{**} , is called the *double dual* or *second dual* of V .

Note that for a finite dimensional space V , $\dim V = \dim V^*$ and also $\dim V^* = \dim V^{**}$, hence V and V^{**} are isomorphic vector spaces. For infinite dimensional spaces $\dim V < \dim V^{**}$ (cf. the exercises) so V and V^{**} cannot be isomorphic. In the case of finite dimensional spaces there is a *natural*, i.e., basis independent or coordinate free way of exhibiting the isomorphism between a vector space and its second dual. The basic idea, in a more general setting, is as follows: if X is any set and S is any set of functions of X into the field F , we normally think of choosing or fixing an $f \in S$ and computing $f(x)$ as x ranges over all of X . Alternatively, we could think of fixing a point x in X and computing $f(x)$ as f ranges over all of S . The latter process, called *evaluation at x* shows that for each $x \in X$ there is a function $E_x : S \rightarrow F$ defined by

$E_x(f) = f(x)$ (i.e., evaluate f at x). This gives a map $x \mapsto E_x$ of X into the set of F -valued functions on S . If S “separates points” in the sense that for distinct points x and y of X there is some $f \in S$ such that $f(x) \neq f(y)$, then the map $x \mapsto E_x$ is injective. The proof of the next lemma applies this “role reversal” process to the situation where $X = V$ and $S = V^*$, proves E_x is a linear F -valued function on S , that is, E_x belongs to the dual space of V^* , and proves the map $x \mapsto E_x$ is a linear transformation from V into V^{**} . Note that throughout this process there is no mention of the word “basis” (although it is convenient to know the dimension of V^{**} — a fact we established by picking bases). In particular, the proof does not start with the familiar phrase “pick a basis of V . . .”

Theorem 19. There is a natural injective linear transformation from V to V^{**} . If V is finite dimensional then this linear transformation is an isomorphism.

Proof: Let $v \in V$. Define the map (*evaluation at v*)

$$E_v : V^* \rightarrow F \quad \text{by} \quad E_v(f) = f(v).$$

Then $E_v(f + \alpha g) = (f + \alpha g)(v) = f(v) + \alpha g(v) = E_v(f) + \alpha E_g(v)$, so that E_v is a linear transformation from V^* to F . Hence E_v is an element of $\text{Hom}_F(V^*, F) = V^{**}$. This defines a natural map

$$\varphi : V \rightarrow V^{**} \quad \text{by} \quad \varphi(v) = E_v.$$

The map φ is a *linear* map, as follows: for $v, w \in V$ and $\alpha \in F$,

$$E_{v+\alpha w}(f) = f(v + \alpha w) = f(v) + \alpha f(w) = E_v(f) + \alpha E_w(f)$$

for every $f \in V^*$, and so

$$\varphi(v + \alpha w) = E_{v+\alpha w} = E_v + \alpha E_w = \varphi(v) + \alpha \varphi(w).$$

To see that φ is injective let v be any nonzero vector in V . By the Building Up Lemma there is a basis \mathcal{B} containing v . Let f be the linear transformation from V to F defined by sending v to 1 and every element of $\mathcal{B} - \{v\}$ to zero. Then $f \in V^*$ and $E_v(f) = f(v) = 1$. Thus $\varphi(v) = E_v$ is not zero in V^{**} . This proves $\ker \varphi = 0$, i.e., φ is injective.

If V has finite dimension n then by Proposition 18, V^* and hence also V^{**} has dimension n . In this case φ is an injective linear transformation from V to a finite dimensional vector space of the same dimension, hence is an isomorphism.

Let V, W be finite dimensional vector spaces over F with bases \mathcal{B}, \mathcal{E} , respectively and let $\mathcal{B}^*, \mathcal{E}^*$ be the dual bases. Fix some $\varphi \in \text{Hom}_F(V, W)$. Then for each $f \in W^*$, the composite $f \circ \varphi$ is a linear transformation from V to F , that is $f \circ \varphi \in V^*$. Thus the map $f \mapsto f \circ \varphi$ defines a function from W^* to V^* . We denote this induced function on dual spaces by φ^* .

Theorem 20. With notations as above, φ^* is a linear transformation from W^* to V^* and $M_{\mathcal{E}^*}^{\mathcal{B}^*}(\varphi^*)$ is the transpose of the matrix $M_{\mathcal{B}}^{\mathcal{E}}(\varphi)$ (recall that the transpose of the matrix (a_{ij}) is the matrix (a_{ji})).

Proof: The map φ^* is linear because $(f + \alpha g) \circ \varphi = (f \circ \varphi) + \alpha(g \circ \varphi)$. The equations which define φ are (from its matrix)

$$\varphi(v_j) = \sum_{i=1}^m \alpha_{ij} w_i \quad 1 \leq j \leq n.$$

To compute the matrix for φ^* , observe that by the definitions of φ^* and w_k^*

$$\varphi^*(w_k^*)(v_j) = (w_k^* \circ \varphi)(v_j) = w_k^* \left(\sum_{i=1}^m \alpha_{ij} w_i \right) = \alpha_{kj}.$$

Also

$$\left(\sum_{i=1}^n \alpha_{ki} v_i^* \right)(v_j) = \alpha_{kj}$$

for all j . This shows that the two linear functionals below agree on a basis of V , hence they are the same element of V^* :

$$\varphi^*(w_k^*) = \sum_{i=1}^n \alpha_{ki} v_i^*.$$

This determines the matrix for φ^* with respect to the bases \mathcal{E}^* and \mathcal{B}^* as the transpose of the matrix for φ .

Corollary 21. For any matrix A , the row rank of A equals the column rank of A .

Proof: Let $\varphi : V \rightarrow W$ be a linear transformation whose matrix with respect to some fixed bases of V and W is A . By Theorem 20 the matrix of $\varphi^* : W^* \rightarrow V^*$ with respect to the dual bases is the transpose of A . The column rank of A is the rank of φ and the row rank of A (= the column rank of the transpose of A) is the rank of φ^* (cf. Exercise 6 of Section 2). It therefore suffices to show that φ and φ^* have the same rank. Now

$$\begin{aligned} f \in \ker \varphi^* &\Leftrightarrow \varphi^*(f) = 0 \Leftrightarrow f \circ \varphi(v) = 0, \quad \text{for all } v \in V \\ &\Leftrightarrow \varphi(V) \subseteq \ker f \Leftrightarrow f \in \text{Ann}(\varphi(V)), \end{aligned}$$

where $\text{Ann}(S)$ is the annihilator of S described in Exercise 3 below. Thus $\text{Ann}(\varphi(V)) = \ker \varphi^*$. By Exercise 3, $\dim \text{Ann}(\varphi(V)) = \dim W - \dim \varphi(V)$. By Corollary 8, $\dim \ker \varphi^* = \dim W^* - \dim \varphi^*(W^*)$. Since W and W^* have the same dimension, $\dim \varphi(V) = \dim \varphi^*(W^*)$ as needed.

EXERCISES

1. Let V be a finite dimensional vector space. Prove that the map $\varphi \mapsto \varphi^*$ in Theorem 20 gives a ring isomorphism of $\text{End}(V)$ with $\text{End}(V^*)$.
2. Let V be the collection of polynomials with coefficients in \mathbb{Q} in the variable x of degree at most 5 with $1, x, x^2, \dots, x^5$ as basis. Prove that the following are elements of the dual space of V and express them as linear combinations of the dual basis:
 - (a) $E : V \rightarrow \mathbb{Q}$ defined by $E(p(x)) = p(3)$ (i.e., evaluation at $x = 3$).
 - (b) $\varphi : V \rightarrow \mathbb{Q}$ defined by $\varphi(p(x)) = \int_0^1 p(t)dt$.
 - (c) $\varphi : V \rightarrow \mathbb{Q}$ defined by $\varphi(p(x)) = \int_0^1 t^2 p(t)dt$.
 - (d) $\varphi : V \rightarrow \mathbb{Q}$ defined by $\varphi(p(x)) = p'(5)$ where $p'(x)$ denotes the usual derivative of the polynomial $p(x)$ with respect to x .
3. Let S be any subset of V^* for some finite dimensional space V . Define $\text{Ann}(S) = \{v \in V \mid f(v) = 0 \text{ for all } f \in S\}$. ($\text{Ann}(S)$ is called the *annihilator of S in V*).
 - (a) Prove that $\text{Ann}(S)$ is a subspace of V .
 - (b) Let W_1 and W_2 be subspaces of V^* . Prove that $\text{Ann}(W_1 + W_2) = \text{Ann}(W_1) \cap \text{Ann}(W_2)$ and $\text{Ann}(W_1 \cap W_2) = \text{Ann}(W_1) + \text{Ann}(W_2)$.
 - (c) Let W_1 and W_2 be subspaces of V^* . Prove that $W_1 = W_2$ if and only if $\text{Ann}(W_1) = \text{Ann}(W_2)$.
 - (d) Prove that the annihilator of S is the same as the annihilator of the subspace of V^* spanned by S .
 - (e) Assume V is finite dimensional with basis v_1, \dots, v_n . Prove that if $S = \{v_1^*, \dots, v_k^*\}$ for some $k \leq n$, then $\text{Ann}(S)$ is the subspace spanned by $\{v_{k+1}, \dots, v_n\}$.
 - (f) Assume V is finite dimensional. Prove that if W^* is any subspace of V^* then $\dim \text{Ann}(W^*) = \dim V - \dim W^*$.
4. If V is infinite dimensional with basis \mathcal{A} , prove that $\mathcal{A}^* = \{v^* \mid v \in \mathcal{A}\}$ does *not* span V^* .
5. If V is infinite dimensional with basis \mathcal{A} , prove that V^* is isomorphic to the direct product of copies of F indexed by \mathcal{A} . Deduce that $\dim V^* > \dim V$. [Use Exercise 14, Section 1.]

11.4 DETERMINANTS

Although we shall be using the theory primarily for vector spaces over a field, the theory of determinants can be developed with no extra effort over arbitrary commutative rings with 1. Thus in this section R is any commutative ring with 1 and V_1, V_2, \dots, V_n, V and W are R -modules. For convenience we repeat the definition of multilinear functions from Section 10.4.

Definition.

- (1) A map $\varphi : V_1 \times V_2 \times \cdots \times V_n \rightarrow W$ is called *multilinear* if for each fixed i and fixed elements $v_j \in V_j, j \neq i$, the map

$$V_i \rightarrow W \quad \text{defined by} \quad x \mapsto \varphi(v_1, \dots, v_{i-1}, x, v_{i+1}, \dots, v_n)$$

is an R -module homomorphism. If $V_i = V, i = 1, 2, \dots, n$, then φ is called an *n -multilinear function on V* , and if in addition $W = R$, φ is called an *n -multilinear form on V* .

- (2) An n -multilinear function φ on V is called *alternating* if $\varphi(v_1, v_2, \dots, v_n) = 0$ whenever $v_i = v_{i+1}$ for some $i \in \{1, 2, \dots, n-1\}$ (i.e., φ is zero whenever two consecutive arguments are equal). The function φ is called *symmetric* if interchanging v_i and v_j for any i and j in (v_1, v_2, \dots, v_n) does not alter the value of φ on this n -tuple.

When $n = 2$ (respectively, 3) one says φ is *bilinear* (respectively, *trilinear*) rather than 2-multilinear (respectively, 3-multilinear). Also, when n is clear from the context we shall simply say φ is multilinear.

Example

For any fixed $m \geq 0$ the usual dot product on $V = \mathbb{R}^m$ is a bilinear form (here the ring R is the field of real numbers).

Proposition 22. Let φ be an n -multilinear alternating function on V . Then

- (1) $\varphi(v_1, \dots, v_{i-1}, v_{i+1}, v_i, v_{i+2}, \dots, v_n) = -\varphi(v_1, v_2, \dots, v_n)$ for any $i \in \{1, 2, \dots, n-1\}$, i.e., the value of φ on an n -tuple is negated if two adjacent components are interchanged.
- (2) For each $\sigma \in S_n$, $\varphi(v_{\sigma(1)}, v_{\sigma(2)}, \dots, v_{\sigma(n)}) = \epsilon(\sigma)\varphi(v_1, v_2, \dots, v_n)$, where $\epsilon(\sigma)$ is the sign of the permutation σ (cf. Section 3.5).
- (3) If $v_i = v_j$ for any pair of distinct $i, j \in \{1, 2, \dots, n\}$ then $\varphi(v_1, v_2, \dots, v_n) = 0$.
- (4) If v_i is replaced by $v_i + \alpha v_j$ in (v_1, \dots, v_n) for any $j \neq i$ and any $\alpha \in R$, the value of φ on this n -tuple is not changed.

Proof: (1) Let $\psi(x, y)$ be the function φ with variable entries x and y in positions i and $i+1$ respectively and fixed entries v_j in position j , for all other j . Thus (1) is the same as showing $\psi(y, x) = -\psi(x, y)$. Since φ is alternating $\psi(x+y, x+y) = 0$. Expanding $x+y$ in each variable in turn gives $\psi(x+y, x+y) = \psi(x, x) + \psi(x, y) + \psi(y, x) + \psi(y, y)$. Again, by the alternating property of φ , the first and last terms on the right hand side of the latter equation are zero. Thus $0 = \psi(x, y) + \psi(y, x)$, which gives (1).

(2) Every permutation can be written as a product of transpositions (cf. Section 3.5). Furthermore, every transposition may be written as a product of transpositions which interchange two successive integers (cf. Exercise 3 of Section 3.5). Thus every permutation σ can be written as $\tau_1 \cdots \tau_m$, where τ_k is a transposition interchanging two successive integers, for all k . It follows from m applications of (1) that

$$\varphi(v_{\sigma(1)}, v_{\sigma(2)}, \dots, v_{\sigma(n)}) = \epsilon(\tau_m) \cdots \epsilon(\tau_1) \varphi(v_1, v_2, \dots, v_n).$$

Finally, since ϵ is a homomorphism into the abelian group ± 1 (so the order of the factors ± 1 does not matter), $\epsilon(\tau_1) \cdots \epsilon(\tau_m) = \epsilon(\tau_1 \cdots \tau_m) = \epsilon(\sigma)$. This proves (2).

(3) Choose σ to be any permutation which fixes i and moves j to $i+1$. Thus $(v_{\sigma(1)}, v_{\sigma(2)}, \dots, v_{\sigma(n)})$ has two equal adjacent components so φ is zero on this n -tuple. By (2), $\varphi(v_{\sigma(1)}, v_{\sigma(2)}, \dots, v_{\sigma(n)}) = \pm \varphi(v_1, v_2, \dots, v_n)$. This implies (3).

(4) This follows immediately from (3) on expanding by linearity in the i^{th} position.

Proposition 23. Assume φ is an n -multilinear alternating function on V and that for some v_1, v_2, \dots, v_n and $w_1, w_2, \dots, w_n \in V$ and some $\alpha_{ij} \in R$ we have

$$\begin{aligned} w_1 &= \alpha_{11}v_1 + \alpha_{21}v_2 + \cdots + \alpha_{n1}v_n \\ w_2 &= \alpha_{12}v_1 + \alpha_{22}v_2 + \cdots + \alpha_{n2}v_n \\ &\vdots \\ w_n &= \alpha_{1n}v_1 + \alpha_{2n}v_2 + \cdots + \alpha_{nn}v_n \end{aligned}$$

(we have purposely written the indices of the α_{ij} in “column format”). Then

$$\varphi(w_1, w_2, \dots, w_n) = \sum_{\sigma \in S_n} \epsilon(\sigma) \alpha_{\sigma(1)1} \alpha_{\sigma(2)2} \cdots \alpha_{\sigma(n)n} \varphi(v_1, v_2, \dots, v_n).$$

Proof: If we expand $\varphi(w_1, w_2, \dots, w_n)$ by multilinearity we obtain a sum of n^n terms of the form $\alpha_{i_11}\alpha_{i_22}\dots\alpha_{i_nn}\varphi(v_{i_1}, v_{i_2}, \dots, v_{i_n})$, where the indices i_1, i_2, \dots, i_n each run over $1, 2, \dots, n$. By Proposition 22(3), φ is zero on the terms where two or more of the i_j 's are equal. Thus in this expansion we need only consider the terms where i_1, \dots, i_n are distinct. Such sequences are in bijective correspondence with permutations in S_n , so each nonzero term may be written as $\alpha_{\sigma(1)1}\alpha_{\sigma(2)2}\cdots\alpha_{\sigma(n)n}\varphi(v_{\sigma(1)}, v_{\sigma(2)}, \dots, v_{\sigma(n)})$, for some $\sigma \in S_n$. Applying (2) of the previous proposition to each of these terms in the expansion of $\varphi(w_1, w_2, \dots, w_n)$ gives the expression in the proposition.

Definition. An $n \times n$ determinant function on R is any function

$$\det : M_{n \times n}(R) \rightarrow R$$

that satisfies the following two axioms:

- (1) \det is an n -multilinear alternating form on $R^n (= V)$, where the n -tuples are the n columns of the matrices in $M_{n \times n}(R)$
- (2) $\det(I) = 1$, where I is the $n \times n$ identity matrix.

On occasion we shall write $\det(A_1, A_2, \dots, A_n)$ for $\det A$, where A_1, A_2, \dots, A_n are the columns of A .

Theorem 24. There is a unique $n \times n$ determinant function on R and it can be computed for any $n \times n$ matrix (α_{ij}) by the formula:

$$\det(\alpha_{ij}) = \sum_{\sigma \in S_n} \epsilon(\sigma) \alpha_{\sigma(1)1} \alpha_{\sigma(2)2} \cdots \alpha_{\sigma(n)n}.$$

Proof: Let A_1, A_2, \dots, A_n be the column vectors in a general $n \times n$ matrix (α_{ij}) . We leave it as an exercise to check that the formula given in the statement of the theorem does satisfy the axioms of a determinant function — this gives existence of a determinant

function. To prove uniqueness let e_i be the column n -tuple with 1 in position i and zeros in all other positions. Then

$$\begin{aligned} A_1 &= \alpha_{11}e_1 + \alpha_{21}e_2 + \cdots + \alpha_{n1}e_n \\ A_2 &= \alpha_{12}e_1 + \alpha_{22}e_2 + \cdots + \alpha_{n2}e_n \\ &\vdots \\ A_n &= \alpha_{1n}e_1 + \alpha_{2n}e_2 + \cdots + \alpha_{nn}e_n. \end{aligned}$$

By Proposition 23, $\det A = \sum_{\sigma \in S_n} \epsilon(\sigma) \alpha_{\sigma(1)1} \alpha_{\sigma(2)2} \cdots \alpha_{\sigma(n)n} \det(e_1, e_2, \dots, e_n)$. Since by axiom (2) of a determinant function $\det(e_1, e_2, \dots, e_n) = 1$, the value of $\det A$ is as claimed.

Corollary 25. The determinant is an n -multilinear function of the rows of $M_{n \times n}(R)$ and for any $n \times n$ matrix A , $\det A = \det(A^t)$, where A^t is the transpose of A .

Proof: The first statement is an immediate consequence of the second, so it suffices to prove that a matrix and its transpose have the same determinant. For $A = (\alpha_{ij})$ one calculates that

$$\det A^t = \sum_{\sigma \in S_n} \epsilon(\sigma) \alpha_{1\sigma(1)} \alpha_{2\sigma(2)} \cdots \alpha_{n\sigma(n)}.$$

Each number from 1 to n appears exactly once among $\sigma(1), \dots, \sigma(n)$ so we may rearrange the product $\alpha_{1\sigma(1)} \alpha_{2\sigma(2)} \cdots \alpha_{n\sigma(n)}$ as $\alpha_{\sigma^{-1}(1)1} \alpha_{\sigma^{-1}(2)2} \cdots \alpha_{\sigma^{-1}(n)n}$. Also, the homomorphism ϵ takes values in $\{\pm 1\}$ so $\epsilon(\sigma) = \epsilon(\sigma^{-1})$. Thus the sum for $\det A^t$ may be rewritten as

$$\sum_{\sigma \in S_n} \epsilon(\sigma^{-1}) \alpha_{\sigma^{-1}(1)1} \alpha_{\sigma^{-1}(2)2} \cdots \alpha_{\sigma^{-1}(n)n}.$$

The latter sum is over all permutations, so the index σ^{-1} may be replaced by σ . The resulting expression is the sum for $\det A$. This completes the proof.

Theorem 26. (Cramer's Rule) If A_1, A_2, \dots, A_n are the columns of an $n \times n$ matrix A and $B = \beta_1 A_1 + \beta_2 A_2 + \cdots + \beta_n A_n$, for some $\beta_1, \dots, \beta_n \in R$, then

$$\beta_i \det A = \det(A_1, \dots, A_{i-1}, B, A_{i+1}, \dots, A_n).$$

Proof: This follows immediately from Proposition 22(3) on replacing the given expression for B in the i^{th} position and expanding by multilinearity in that position.

Corollary 27. If R is an integral domain, then $\det A = 0$ for $A \in M_n(R)$ if and only if the columns of A are R -linearly dependent as elements of the free R -module of rank n . Also, $\det A = 0$ if and only if the rows of A are R -linearly dependent.

Proof: Since $\det A = \det A^t$ the first sentence implies the second.

Assume first that the columns of A are linearly dependent and

$$0 = \beta_1 A_1 + \beta_2 A_2 + \cdots + \beta_n A_n$$

is a dependence relation on the columns of A with, say, $\beta_i \neq 0$. By Cramer's Rule, $\beta_i \det A = 0$. Since R is an integral domain and $\beta_i \neq 0$, $\det A = 0$.

Conversely, assume the columns of A are independent. Consider the integral domain R as embedded in its quotient field F so that $M_{n \times n}(R)$ may be considered as a subring of $M_{n \times n}(F)$ (and note that the determinant function on the subring is the restriction of the determinant function from $M_{n \times n}(F)$). The columns of A in this way become elements of F^n . Any nonzero F -linear combination of the columns of A which is zero in F^n gives, by multiplying the coefficients by a common denominator, a nonzero R -linear dependence relation. The columns of A must therefore be independent vectors in F^n . Since A has n columns, these form a basis of F^n . Thus there are elements β_{ij} of F such that for each i , the i^{th} basis vector e_i in F^n may be expressed as

$$e_i = \beta_{1i}A_1 + \beta_{2i}A_2 + \cdots + \beta_{ni}A_n.$$

The $n \times n$ identity matrix is the one whose columns are e_1, e_2, \dots, e_n . By Proposition 23 (with $\varphi = \det$), the determinant of the identity matrix is some F -multiple of $\det A$. Since the determinant of the identity matrix is 1, $\det A$ cannot be zero. This completes the proof.

Theorem 28. For matrices $A, B \in M_{n \times n}(R)$, $\det AB = (\det A)(\det B)$.

Proof: Let $B = (\beta_{ij})$ and let A_1, A_2, \dots, A_n be the columns of A . Then $C = AB$ is the $n \times n$ matrix whose j^{th} column is $C_j = \beta_{1j}A_1 + \beta_{2j}A_2 + \cdots + \beta_{nj}A_n$. By Proposition 23 applied to the multilinear function \det we obtain

$$\det C = \det(C_1, \dots, C_n) = \left[\sum_{\sigma \in S_n} \epsilon(\sigma) \beta_{\sigma(1)1} \beta_{\sigma(2)2} \cdots \beta_{\sigma(n)n} \right] \det(A_1, \dots, A_n).$$

The sum inside the brackets is the formula for $\det B$, hence $\det C = (\det B)(\det A)$, as required (R is commutative).

Definition. Let $A = (\alpha_{ij})$ be an $n \times n$ matrix. For each i, j , let A_{ij} be the $(n-1) \times (n-1)$ matrix obtained from A by deleting its i^{th} row and j^{th} column (an $(n-1) \times (n-1)$ minor of A). Then $(-1)^{i+j} \det(A_{ij})$ is called the ij cofactor of A .

Theorem 29. (The Cofactor Expansion Formula along the i^{th} row) If $A = (\alpha_{ij})$ is an $n \times n$ matrix, then for each fixed $i \in \{1, 2, \dots, n\}$ the determinant of A can be computed from the formula

$$\det A = (-1)^{i+1} \alpha_{i1} \det A_{i1} + (-1)^{i+2} \alpha_{i2} \det A_{i2} + \cdots + (-1)^{i+n} \alpha_{in} \det A_{in}.$$

Proof: For each A let $D(A)$ be the element of R obtained from the cofactor expansion formula described above. We prove that D satisfies the axioms of a determinant function, hence is the determinant function. Proceed by induction on n . If $n = 1$, $D((\alpha)) = \alpha$, for all 1×1 matrices (α) and the result holds. Assume therefore that $n \geq 2$. To show that D is an alternating multilinear function of the columns, fix an index k and consider the k^{th} column as varying and all other columns as fixed. If $j \neq k$,

α_{ij} does not depend on k and $D(A_{ij})$ is linear in the k^{th} column by induction. Also, as the k^{th} column varies linearly so does α_{ik} , whereas $D(A_{ik})$ remains unchanged (the k^{th} column has been deleted from A_{ik}). Thus each term in the formula for D varies linearly in the k^{th} column. This proves D is multilinear in the columns.

To prove D is alternating assume columns k and $k+1$ of A are equal. If $j \neq k$ or $k+1$, the two equal columns of A become two equal columns in the matrix A_{ij} . By induction $D(A_{ij}) = 0$. The formula for D therefore has at most two nonzero terms: when $j = k$ and when $j = k+1$. The minor matrices A_{ik} and $A_{i,k+1}$ are identical and $\alpha_{ik} = \alpha_{i,k+1}$. Then the two remaining terms in the expansion for D , $(-1)^{i+k}\alpha_{ik}D(A_{ik})$ and $(-1)^{i+k+1}\alpha_{i,k+1}D(A_{i,k+1})$ are equal and appear with opposite signs, hence they cancel. Thus $D(A) = 0$ if A has two adjacent columns which are equal, i.e., D is alternating.

Finally, it follows easily from the formula and induction that $D(I) = 1$, where I is the identity matrix. This completes the induction.

Theorem 30. (Cofactor Formula for the Inverse of a Matrix) Let $A = (\alpha_{ij})$ be an $n \times n$ matrix and let B be the transpose of its matrix of cofactors, i.e., $B = (\beta_{ij})$, where $\beta_{ij} = (-1)^{i+j} \det A_{ji}$, $1 \leq i, j \leq n$. Then $AB = BA = (\det A)I$. Moreover, $\det A$ is a unit in R if and only if A is a unit in $M_{n \times n}(R)$; in this case the matrix $\frac{1}{\det A}B$ is the inverse of A .

Proof: The i, j entry of AB is $\alpha_{i1}\beta_{1j} + \alpha_{i2}\beta_{2j} + \cdots + \alpha_{in}\beta_{nj}$. By definition of the entries of B this equals

$$\alpha_{i1}(-1)^{j+1}D(A_{j1}) + \alpha_{i2}(-1)^{j+2}D(A_{j2}) + \cdots + \alpha_{in}(-1)^{j+n}D(A_{jn}). \quad (11.7)$$

If $i = j$, this is the cofactor expansion for $\det A$ along the i^{th} row. The diagonal entries of AB are thus all equal to $\det A$. If $i \neq j$, let \bar{A} be the matrix A with the j^{th} row replaced by the i^{th} row, so $\det \bar{A} = 0$. By inspection $\bar{A}_{jk} = A_{jk}$ and $\alpha_{ik} = \bar{\alpha}_{jk}$ for every $k \in \{1, 2, \dots, n\}$. By making these substitutions in equation (7) for each $k = 1, 2, \dots, n$ one sees that the i, j entry in AB equals $\bar{\alpha}_{j1}(-1)^{1+j}D(\bar{A}_{j1}) + \cdots + \bar{\alpha}_{jn}(-1)^{n+j}D(\bar{A}_{jn})$. This expression is the cofactor expansion for $\det \bar{A}$ along the j^{th} row. Since, as noted above, $\det \bar{A} = 0$, this proves that all off diagonal terms of AB are zero, which proves that $AB = (\det A)I$.

It follows directly from the definition of B that the pair (A', B') satisfies the same hypotheses as the pair (A, B) . By what has already been shown it follows that $(BA)' = A'B' = (\det A')I$. Since $\det A' = \det A$ and the transpose of a diagonal matrix is itself, we obtain $BA = (\det A)I$ as well.

If $d = \det A$ is a unit in R , then $d^{-1}B$ is a matrix with entries in R whose product with A (on either side) is the identity, i.e., A is a unit in $M_{n \times n}(R)$. Conversely, assume that A is a unit in R with (2-sided) inverse matrix C . Since $\det C \in R$ and

$$1 = \det I = \det AC = (\det A)(\det C) = (\det C)(\det A),$$

it follows that $\det A$ has a 2-sided inverse in R , as needed. This completes all parts of the proof.

EXERCISES

1. Formulate and prove the cofactor expansion formula along the j^{th} column of a square matrix A .
2. Let F be a field and let A_1, A_2, \dots, A_n be (column) vectors in F^n . Form the matrix A whose i^{th} column is A_i . Prove that these vectors form a basis of F^n if and only if $\det A \neq 0$.
3. Let R be any commutative ring with 1, let V be an R -module and let $x_1, x_2, \dots, x_n \in V$. Assume that for some $A \in M_{n \times n}(R)$,

$$A \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = 0.$$

Prove that $(\det A)x_i = 0$, for all $i \in \{1, 2, \dots, n\}$.

4. (*Computing Determinants of Matrices*) This exercise outlines the use of Gauss–Jordan elimination (cf. the exercises in Section 2) to compute determinants. This is the most efficient general procedure for computing large determinants. Let A be an $n \times n$ matrix.
 - (a) Prove that the elementary row operations have the following effect on determinants:
 - (i) interchanging two rows changes the sign of the determinant
 - (ii) adding a multiple of one row to another does not alter the determinant
 - (iii) multiplying any row by a nonzero element u from F multiplies the determinant by u .
 - (b) Prove that $\det A$ is nonzero if and only if A is row equivalent to the $n \times n$ identity matrix. Suppose A can be row reduced to the identity matrix using a total of s row interchanges as in (i) and by multiplying rows by the nonzero elements u_1, u_2, \dots, u_t as in (iii). Prove that $\det A = (-1)^s(u_1u_2 \dots u_t)^{-1}$.
5. Compute the determinants of the following matrices using row reduction:

$$A = \begin{pmatrix} 5 & 4 & -6 \\ -2 & 0 & 2 \\ 3 & 4 & -2 \end{pmatrix} \quad B = \begin{pmatrix} 1 & 2 & -4 & 4 \\ 2 & -1 & 4 & -8 \\ 1 & 0 & 1 & -2 \\ 0 & 1 & -2 & 3 \end{pmatrix}.$$

6. (*Minkowski's Criterion*) Suppose A is an $n \times n$ matrix with real entries such that the diagonal elements are all positive, the off-diagonal elements are all negative and the row sums are all positive. Prove that $\det A \neq 0$. [Consider the corresponding system of equations $AX = 0$ and suppose there is a nontrivial solution (x_1, \dots, x_n) . If x_i has the largest absolute value show that the i^{th} equation leads to a contradiction.]

11.5 TENSOR ALGEBRAS, SYMMETRIC AND EXTERIOR ALGEBRAS

In this section R is any commutative ring with 1, and we assume the left and right actions of R on each R -module are the same. We shall primarily be interested in the special case when $R = F$ is a field, but the basic constructions hold in general.

Suppose M is an R -module. When tensor products were first introduced in Section 10.4 we spoke heuristically of forming “products” m_1m_2 of elements of M , and we constructed a new module $M \otimes M$ generated by such “products” $m_1 \otimes m_2$. The “value” of this product is not in M , so this does not give a ring structure on M itself. If, however,

we iterate this by taking the “products” $m_1 m_2 m_3$ and $m_1 m_2 m_3 m_4$, and all finite sums of such products, we can construct a ring containing M that is “universal” with respect to rings containing M (and, more generally, with respect to homomorphic images of M), as we now show.

For each integer $k \geq 1$, define

$$\mathcal{T}^k(M) = M \otimes_R M \otimes_R \cdots \otimes_R M \quad (k \text{ factors}),$$

and set $\mathcal{T}^0(M) = R$. The elements of $\mathcal{T}^k(M)$ are called k -tensors. Define

$$\mathcal{T}(M) = R \oplus \mathcal{T}^1(M) \oplus \mathcal{T}^2(M) \oplus \mathcal{T}^3(M) \cdots = \bigoplus_{k=0}^{\infty} \mathcal{T}^k(M).$$

Every element of $\mathcal{T}(M)$ is a finite linear combination of k -tensors for various $k \geq 0$. We identify M with $\mathcal{T}^1(M)$, so that M is an R -submodule of $\mathcal{T}(M)$.

Theorem 31. If M is any R -module over the commutative ring R then

(1) $\mathcal{T}(M)$ is an R -algebra containing M with multiplication defined by mapping

$$(m_1 \otimes \cdots \otimes m_i)(m'_1 \otimes \cdots \otimes m'_j) = m_1 \otimes \cdots \otimes m_i \otimes m'_1 \otimes \cdots \otimes m'_j$$

and extended to sums via the distributive laws. With respect to this multiplication $\mathcal{T}^i(M)\mathcal{T}^j(M) \subseteq \mathcal{T}^{i+j}(M)$.

(2) (*Universal Property*) If A is any R -algebra and $\varphi : M \rightarrow A$ is an R -module homomorphism, then there is a unique R -algebra homomorphism $\Phi : \mathcal{T}(M) \rightarrow A$ such that $\Phi|_M = \varphi$.

Proof: The map

$$\underbrace{M \times M \times \cdots \times M}_{i \text{ factors}} \times \underbrace{M \times M \times \cdots \times M}_{j \text{ factors}} \rightarrow \mathcal{T}^{i+j}(M)$$

defined by

$$(m_1, \dots, m_i, m'_1, \dots, m'_j) \mapsto m_1 \otimes \cdots \otimes m_i \otimes m'_1 \otimes \cdots \otimes m'_j$$

is R -multilinear, so induces a bilinear map $\mathcal{T}^i(M) \times \mathcal{T}^j(M)$ to $\mathcal{T}^{i+j}(M)$ which is easily checked to give a well defined multiplication satisfying (1) (cf. the proof of Proposition 21 in Section 10.4). To prove (2), assume that $\varphi : M \rightarrow A$ is an R -algebra homomorphism. Then

$$(m_1, m_2, \dots, m_k) \mapsto \varphi(m_1)\varphi(m_2)\dots\varphi(m_k)$$

defines an R -multilinear map from $M \times \cdots \times M$ (k times) to A . This in turn induces a unique R -module homomorphism Φ from $\mathcal{T}^k(M)$ to A (Corollary 16 of Section 10.4) mapping $m_1 \otimes \cdots \otimes m_k$ to the element on the right hand side above. It is easy to check from the definition of the multiplication in (1) that the resulting uniquely defined map $\Phi : \mathcal{T}(M) \rightarrow A$ is an R -algebra homomorphism.

Definition. The ring $\mathcal{T}(M)$ is called the *tensor algebra* of M .

Proposition 32. Let V be a finite dimensional vector space over the field F with basis $\mathcal{B} = \{v_1, \dots, v_n\}$. Then the k -tensors

$$v_{i_1} \otimes v_{i_2} \otimes \cdots \otimes v_{i_k} \quad \text{with } v_{i_j} \in \mathcal{B}$$

are a vector space basis of $\mathcal{T}^k(V)$ over F (with the understanding that the basis vector is the element $1 \in F$ when $k = 0$). In particular, $\dim_F(\mathcal{T}^k(V)) = n^k$.

Proof: This follows immediately from Proposition 16 of Section 2.

Theorem 31 and Proposition 32 show that the space $\mathcal{T}(V)$ may be regarded as the *noncommutative polynomial algebra* over F in the (noncommuting) variables v_1, \dots, v_n . The analogous result also holds for finitely generated free modules over any commutative ring (using Corollary 19 in Section 10.4).

Examples

- (1) Let $R = \mathbb{Z}$ and let $M = \mathbb{Q}/\mathbb{Z}$. Then $(\mathbb{Q}/\mathbb{Z}) \otimes_{\mathbb{Z}} (\mathbb{Q}/\mathbb{Z}) = 0$ (Example 4 following Corollary 12 in Section 10.4). Thus $\mathcal{T}(\mathbb{Q}/\mathbb{Z}) = \mathbb{Z} \oplus (\mathbb{Q}/\mathbb{Z})$, where addition is componentwise and the multiplication is given by $(r, \bar{p})(s, \bar{q}) = (rs, \overline{rq + sp})$. The ring $R/(x)$ of Exercise 4(d) in Section 9.3 is isomorphic to $\mathcal{T}(\mathbb{Q}/\mathbb{Z})$.
- (2) Let $R = \mathbb{Z}$ and let $M = \mathbb{Z}/n\mathbb{Z}$. Then $(\mathbb{Z}/n\mathbb{Z}) \otimes_{\mathbb{Z}} (\mathbb{Z}/n\mathbb{Z}) \cong \mathbb{Z}/n\mathbb{Z}$ (Example 3 following Corollary 12 in Section 10.4). Thus $\mathcal{T}^i(M) \cong M$ for all $i > 0$ and so $\mathcal{T}(\mathbb{Z}/n\mathbb{Z}) \cong \mathbb{Z} \oplus (\mathbb{Z}/n\mathbb{Z}) \oplus (\mathbb{Z}/n\mathbb{Z}) \cdots$. It follows easily that $\mathcal{T}(\mathbb{Z}/n\mathbb{Z}) \cong \mathbb{Z}[x]/(nx)$.

Since $\mathcal{T}^i(M)\mathcal{T}^j(M) \subseteq \mathcal{T}^{i+j}(M)$, the tensor algebra $\mathcal{T}(M)$ has a natural “grading” or “degree” structure reminiscent of a polynomial ring.

Definition.

- (1) A ring S is called a *graded ring* if it is the direct sum of additive subgroups: $S = S_0 \oplus S_1 \oplus S_2 \oplus \cdots$ such that $S_i S_j \subseteq S_{i+j}$ for all $i, j \geq 0$. The elements of S_k are said to be *homogeneous of degree k* , and S_k is called the *homogeneous component of S of degree k* .
- (2) An ideal I of the graded ring S is called a *graded ideal* if $I = \bigoplus_{k=0}^{\infty} (I \cap S_k)$.
- (3) A ring homomorphism $\varphi : S \rightarrow T$ between two graded rings is called a *homomorphism of graded rings* if it respects the grading structures on S and T , i.e., if $\varphi(S_k) \subseteq T_k$ for $k = 0, 1, 2, \dots$.

Note that $S_0 S_0 \subseteq S_0$, which implies that S_0 is a subring of the graded ring S and then S is an S_0 -module. If S_0 is in the center of S and it contains an identity of S , then S is an S_0 -algebra. Note also that the ideal I is graded if whenever a sum $i_{k_1} + \cdots + i_{k_n}$ of homogeneous elements with distinct degrees k_1, \dots, k_n is in I then each of the individual summands i_{k_1}, \dots, i_{k_n} is itself in I .

Example

The polynomial ring $S = R[x_1, x_2, \dots, x_n]$ in n variables over the commutative ring R is an example of a graded ring. Here $S_0 = R$ and the homogeneous component of degree k is the subgroup of all R -linear combinations of monomials of degree k .

The ideal I generated by x_1, \dots, x_n is a graded ideal: every polynomial with zero constant term may be written uniquely as a sum of homogeneous polynomials of degree $k > 1$, and each of these has zero constant term hence lies in I . More generally, an ideal is a graded ideal if and only if it can be generated by homogeneous polynomials (cf. Exercise 17 in Section 9.1).

Not every ideal of a graded ring need be a graded ideal. For example in the graded ring $\mathbb{Z}[x]$ the principal ideal J generated by $1 + x$ is not graded: $1 + x \in J$ and $1 \notin J$ so $1 + x$ cannot be written as a sum of homogeneous polynomials each of which belongs to J .

The next result shows that quotients of graded rings by graded ideals are again graded rings.

Proposition 33. Let S be a graded ring, let I be a graded ideal in S and let $I_k = I \cap S_k$ for all $k \geq 0$. Then S/I is naturally a graded ring whose homogeneous component of degree k is isomorphic to S_k/I_k .

Proof: The map

$$\begin{aligned} S &= \bigoplus_{k=0}^{\infty} S_k \longrightarrow \bigoplus_{k=0}^{\infty} (S_k/I_k) \\ (\dots, s_k, \dots) &\longmapsto (\dots, s_k \bmod I_k, \dots) \end{aligned}$$

is surjective with kernel $I = \bigoplus_{k=0}^{\infty} I_k$ and defines an isomorphism of graded rings. The details are left for the exercises.

Symmetric Algebras

The first application of Proposition 33 is in the construction of a commutative quotient ring of $\mathcal{T}(M)$ through which R -module homomorphisms from M to any *commutative* R -algebra must factor. This gives an “abelianized” version of Theorem 31. The construction is analogous to forming the commutator quotient G/G' of a group (cf. Section 5.4).

Definition. The *symmetric algebra* of an R -module M is the R -algebra obtained by taking the quotient of the tensor algebra $\mathcal{T}(M)$ by the ideal $\mathcal{C}(M)$ generated by all elements of the form $m_1 \otimes m_2 - m_2 \otimes m_1$, for all $m_1, m_2 \in M$. The symmetric algebra $\mathcal{T}(M)/\mathcal{C}(M)$ is denoted by $\mathcal{S}(M)$.

The tensor algebra $\mathcal{T}(M)$ is generated as a ring by $R = \mathcal{T}^0(M)$ and $M = \mathcal{T}^1(M)$, and these elements commute in the quotient ring $\mathcal{S}(M)$ by definition. It follows that the symmetric algebra $\mathcal{S}(M)$ is a commutative ring. The ideal $\mathcal{C}(M)$ is generated by homogeneous tensors of degree 2 and it follows easily that $\mathcal{C}(M)$ is a graded ideal. Then by Proposition 33 the symmetric algebra is a graded ring whose homogeneous component of degree k is $\mathcal{S}^k(M) = \mathcal{T}^k(M)/\mathcal{C}^k(M)$. Since $\mathcal{C}(M)$ consists of k -tensors

with $k \geq 2$, we have $\mathcal{C}(M) \cap M = 0$ and so the image of $M = \mathcal{T}^1(M)$ in $\mathcal{S}(M)$ is isomorphic to M . Identifying M with its image we see that $\mathcal{S}^1(M) = M$ and the symmetric algebra contains M . In a similar way $\mathcal{S}^0(M) = R$, so the symmetric algebra is also an R -algebra. The R -module $\mathcal{S}^k(M)$ is called the k^{th} symmetric power of M .

The first part of the next theorem shows that the elements of the k^{th} symmetric power of M can be considered as finite sums of simple tensors $m_1 \otimes \cdots \otimes m_k$ where tensors with the order of the factors permuted are identified. Recall also from Section 4 that a k -multilinear map $\varphi : M \times \cdots \times M \rightarrow N$ is said to be *symmetric* if $\varphi(m_1, \dots, m_k) = \varphi(m_{\sigma(1)}, \dots, m_{\sigma(k)})$ for all permutations σ of $1, 2, \dots, k$. (The definition is the same for modules over any commutative ring R as for vector spaces.)

Theorem 34. Let M be an R -module over the commutative ring R and let $\mathcal{S}(M)$ be its symmetric algebra.

- (1) The k^{th} symmetric power, $\mathcal{S}^k(M)$, of M is equal to $M \otimes \cdots \otimes M$ (k factors) modulo the submodule generated by all elements of the form

$$(m_1 \otimes m_2 \otimes \cdots \otimes m_k) - (m_{\sigma(1)} \otimes m_{\sigma(2)} \otimes \cdots \otimes m_{\sigma(k)})$$

for all $m_i \in M$ and all permutations σ in the symmetric group S_k .

- (2) (Universal Property for Symmetric Multilinear Maps) If $\varphi : M \times \cdots \times M \rightarrow N$ is a symmetric k -multilinear map over R then there is a unique R -module homomorphism $\Phi : \mathcal{S}^k(M) \rightarrow N$ such that $\varphi = \Phi \circ \iota$, where

$$\iota : M \times \cdots \times M \rightarrow \mathcal{S}^k(M)$$

is the map defined by

$$\iota(m_1, \dots, m_k) = m_1 \otimes \cdots \otimes m_n \bmod \mathcal{C}(M).$$

- (3) (Universal Property for maps to commutative R -algebras) If A is any commutative R -algebra and $\varphi : M \rightarrow A$ is an R -module homomorphism, then there is a unique R -algebra homomorphism $\Phi : \mathcal{S}(M) \rightarrow A$ such that $\Phi|_M = \varphi$.

Proof: The k -tensors $\mathcal{C}^k(M)$ in the ideal $\mathcal{C}(M)$ are finite sums of elements of the form

$$m_1 \otimes \cdots \otimes m_{i-1} \otimes (m_i \otimes m_{i+1} - m_{i+1} \otimes m_i) \otimes m_{i+2} \otimes \cdots \otimes m_k$$

with $m_1, \dots, m_k \in M$ (where $k \geq 2$ and $1 \leq i < k$). This product gives a difference of two k -tensors which are equal except that two entries (in positions i and $i+1$) have been transposed, i.e., gives the element in (1) of the theorem corresponding to the transposition $(i \ i+1)$ in the symmetric group S_k . Conversely, since any permutation σ in S_k can be written as a product of such transpositions it is easy to see that every element in (1) can be written as a sum of elements of the form above. This gives (1).

The proofs of (2) and (3) are very similar to the proofs of the corresponding “asymmetric” results (Corollary 16 of Section 10.4 and Theorem 31) noting that $\mathcal{C}^k(M)$ is contained in the kernel of any symmetric map from $\mathcal{T}^k(M)$ to N by part (1).

Corollary 35. Let V be an n -dimensional vector space over the field F . Then $\mathcal{S}(V)$ is isomorphic as a graded F -algebra to the ring of polynomials in n variables over F (i.e., the isomorphism is also a vector space isomorphism from $\mathcal{S}^k(V)$ onto the space of all homogeneous polynomials of degree k). In particular, $\dim_F(\mathcal{S}^k(V)) = \binom{k+n-1}{n-1}$.

Proof: Let $\mathcal{B} = \{v_1, \dots, v_n\}$ be a basis of V . By Proposition 32 there is a bijection between a basis of $\mathcal{T}^k(V)$ and the set \mathcal{B}^k of ordered k -tuples of elements from \mathcal{B} . Define two k -tuples in \mathcal{B}^k to be equivalent if there is some permutation of the entries of one that gives the other — this is easily seen to be an equivalence relation on \mathcal{B}^k . Let $S(\mathcal{B}^k)$ denote the corresponding set of equivalence classes. Any symmetric k -multilinear function from V^k to a vector space over F will be constant on all of the basis tensors whose corresponding k -tuples lie in the same equivalence class; conversely, any function from $S(\mathcal{B}^k)$ can be uniquely extended to a symmetric k -multilinear function on V^k . It follows that the vector space over F with basis $S(\mathcal{B}^k)$ satisfies the universal property of $\mathcal{S}^k(V)$ in Theorem 34(2), hence is isomorphic to $\mathcal{S}^k(V)$. Each equivalence class has a unique representative of the form $(v_1^{a_1}, v_2^{a_2}, \dots, v_n^{a_n})$, where v_i^a denotes the sequence v_i, v_i, \dots, v_i taken a times, each $a_i \geq 0$, and $a_1 + \dots + a_n = k$. Thus there is a bijection between the basis $S^k(\mathcal{B})$ and the set $x_1^{a_1} \cdots x_n^{a_n}$ of monic monomials of degree k in the polynomial ring $F[x_1, \dots, x_n]$. This bijection extends to an isomorphism of graded F -algebras, proving the first part of the corollary. The computation of the dimension of $\mathcal{S}^k(V)$ (i.e., the number of monic monomials of degree k) is left as an exercise.

Exterior Algebras

Recall from Section 4 that a multilinear map $\varphi : M \times \dots \times M \rightarrow N$ is called *alternating* if $\varphi(m_1, \dots, m_k) = 0$ whenever $m_i = m_{i+1}$ for some i . (The definition is the same for any R -module as for vector spaces.) We saw that the determinant map was alternating, and was uniquely determined by some additional constraints. We can apply Proposition 33 to construct an algebra through which alternating multilinear maps must factor in a manner similar to the construction of the symmetric algebra (through which symmetric multilinear maps factor).

Definition. The *exterior algebra* of an R -module M is the R -algebra obtained by taking the quotient of the tensor algebra $\mathcal{T}(M)$ by the ideal $\mathcal{A}(M)$ generated by all elements of the form $m \otimes m$, for $m \in M$. The exterior algebra $\mathcal{T}(M)/\mathcal{A}(M)$ is denoted by $\bigwedge(M)$ and the image of $m_1 \otimes m_2 \otimes \dots \otimes m_k$ in $\bigwedge(M)$ is denoted by $m_1 \wedge m_2 \wedge \dots \wedge m_k$.

As with the symmetric algebra, the ideal $\mathcal{A}(M)$ is generated by homogeneous elements hence is a graded ideal. By Proposition 33 the exterior algebra is graded, with k^{th} homogeneous component $\bigwedge^k(M) = \mathcal{T}^k(M)/\mathcal{A}^k(M)$. We can again identify R with $\bigwedge^0(M)$ and M with $\bigwedge^1(M)$ and so consider M as an R -submodule of the R -algebra $\bigwedge(M)$. The R -module $\bigwedge^k(M)$ is called the k^{th} *exterior power* of M .

The multiplication

$$(m_1 \wedge \dots \wedge m_i) \wedge (m'_1 \wedge \dots \wedge m'_j) = m_1 \wedge \dots \wedge m_i \wedge m'_1 \wedge \dots \wedge m'_j$$

in the exterior algebra is called the *wedge* (or *exterior*) *product*. By definition of the quotient, this multiplication is alternating in the sense that the product $m_1 \wedge \cdots \wedge m_k$ is 0 in $\bigwedge(M)$ if $m_i = m_{i+1}$ for any $1 \leq i < k$. Then

$$\begin{aligned} 0 &= (m + m') \wedge (m + m') \\ &= (m \wedge m) + (m \wedge m') + (m' \wedge m) + (m' \wedge m') \\ &= (m \wedge m') + (m' \wedge m) \end{aligned}$$

shows that the multiplication is also anticommutative on simple tensors:

$$m \wedge m' = -m' \wedge m \quad \text{for all } m, m' \in M.$$

This anticommutativity does not extend to arbitrary products, however, i.e., we need not have $ab = -ba$ for all $a, b \in \bigwedge(M)$ (cf. Exercise 4).

Theorem 36. Let M be an R -module over the commutative ring R and let $\bigwedge(M)$ be its exterior algebra.

- (1) The k^{th} exterior power, $\bigwedge^k(M)$, of M is equal to $M \otimes \cdots \otimes M$ (k factors) modulo the submodule generated by all elements of the form

$$m_1 \otimes m_2 \otimes \cdots \otimes m_k \quad \text{where } m_i = m_j \text{ for some } i \neq j.$$

In particular,

$$m_1 \wedge m_2 \wedge \cdots \wedge m_k = 0 \quad \text{if } m_i = m_j \text{ for some } i \neq j.$$

- (2) (*Universal Property for Alternating Multilinear Maps*) If $\varphi : M \times \cdots \times M \rightarrow N$ is an alternating k -multilinear map then there is a unique R -module homomorphism $\Phi : \bigwedge^k(M) \rightarrow N$ such that $\varphi = \Phi \circ \iota$, where

$$\iota : M \times \cdots \times M \rightarrow \bigwedge^k(M)$$

is the map defined by

$$\iota(m_1, \dots, m_k) = m_1 \wedge \cdots \wedge m_k.$$

Remark: The exterior algebra also satisfies a universal property similar to (3) of Theorem 34, namely with respect to R -module homomorphisms from M to R -algebras A satisfying $a^2 = 0$ for all $a \in A$ (cf. Exercise 6).

Proof: The k -tensors $\mathcal{A}^k(M)$ in the ideal $\mathcal{A}(M)$ are finite sums of elements of the form

$$m_1 \otimes \cdots \otimes m_{i-1} \otimes (m \otimes m) \otimes m_{i+2} \otimes \cdots \otimes m_k$$

with $m_1, \dots, m_k, m \in M$ (where $k \geq 2$ and $1 \leq i < k$), which is a k -tensor with two equal entries (in positions i and $i+1$), so is of the form in (1). For the reverse inclusion, note that since

$$\begin{aligned} m' \otimes m &= -m \otimes m' + [(m + m') \otimes (m + m') - m \otimes m - m' \otimes m'] \\ &\equiv -m \otimes m' \pmod{\mathcal{A}(M)}, \end{aligned}$$

interchanging any two consecutive entries and multiplying by -1 in a simple k -tensor gives an equivalent tensor modulo $\mathcal{A}^k(M)$. Using such a sequence of interchanges and sign changes we can arrange for the equal entries m_i and m_j of a simple tensor as in (1) to be adjacent, which gives an element of $\mathcal{A}^k(M)$. It follows that the generators in (1) are contained in $\mathcal{A}^k(M)$, which proves the first part of the theorem.

As in Theorem 34, the proof of (2) follows easily from the corresponding result for the tensor algebra in Theorem 31 since $\mathcal{A}^k(M)$ is contained in the kernel of any alternating map from $\mathcal{T}^k(M)$ to N .

Examples

- (1) Suppose V is a one-dimensional vector space over F with basis element v . Then $\bigwedge^k(V)$ consists of finite sums of elements of the form $\alpha_1 v \wedge \alpha_2 v \wedge \cdots \wedge \alpha_k v$, i.e., $\alpha_1 \alpha_2 \cdots \alpha_k (v \wedge v \wedge \cdots \wedge v)$ for $\alpha_1, \dots, \alpha_k \in F$. Since $v \wedge v = 0$, it follows that $\bigwedge^0(V) = F$, $\bigwedge^1(V) = V$, and $\bigwedge^i(V) = 0$ for $i \geq 2$, so as a graded F -algebra we have

$$\bigwedge(V) = F \oplus V \oplus 0 \oplus 0 \oplus \dots$$

- (2) Suppose now that V is a two-dimensional vector space over F with basis v, v' . Here $\bigwedge^k(V)$ consists of finite sums of elements of the form $(\alpha_1 v + \alpha'_1 v') \wedge \cdots \wedge (\alpha_k v + \alpha'_k v')$. Such an element is a sum of elements that are simple wedge products involving only v and v' . For example, an element in $\bigwedge^2(V)$ is a sum of elements of the form

$$\begin{aligned} (av + bv') \wedge (cv + dv') &= ac(v \wedge v) + ad(v \wedge v') + bc(v' \wedge v) \\ &\quad + bd(v' \wedge v') \\ &= (ad - bc)v \wedge v'. \end{aligned}$$

It follows that $\bigwedge^i(V) = 0$ for $i \geq 3$ since then at least one of v, v' appears twice in such simple products.

We can see directly from $\bigwedge^2(V) = \mathcal{T}^2(V)/\mathcal{A}^2(V)$ that $v \wedge v' \neq 0$, as follows. The vector space $\mathcal{T}^2(V)$ is 4-dimensional with $v \otimes v, v \otimes v', v' \otimes v, v' \otimes v'$ as basis (Proposition 16). The elements $v \otimes v, v \otimes v' + v' \otimes v, v' \otimes v'$ and $v \otimes v'$ are therefore also a basis for $\mathcal{T}^2(V)$. The subspace $\mathcal{A}^2(V)$ consists of all the 2-tensors in the ideal generated by the tensors

$$(av + bv') \otimes (av + bv') = a^2(v \otimes v) + ab(v \otimes v' + v' \otimes v) + b^2(v' \otimes v'),$$

from which it is clear that $\mathcal{A}^2(V)$ is contained in the 3-dimensional subspace having $v \otimes v, v \otimes v' + v' \otimes v$, and $v' \otimes v'$ as basis. In particular, the basis element $v \otimes v'$ of $\mathcal{T}^2(V)$ is not contained in $\mathcal{A}^2(V)$, i.e., $v \wedge v' \neq 0$ in $\bigwedge^2(V)$.

It follows that $\bigwedge^0(V) = F$, $\bigwedge^1(V) = V$, $\bigwedge^2(V) = F(v \wedge v')$, and $\bigwedge^i(V) = 0$ for $i \geq 3$, so as a graded F -algebra we have

$$\bigwedge(V) = F \oplus V \oplus F(v \wedge v') \oplus 0 \oplus \dots$$

As the previous examples illustrate, unlike the tensor and symmetric algebras, for finite dimensional vector spaces the exterior algebra is finite dimensional:

Corollary 37. Let V be a finite dimensional vector space over the field F with basis $\mathcal{B} = \{v_1, \dots, v_n\}$. Then the vectors

$$v_{i_1} \wedge v_{i_2} \wedge \cdots \wedge v_{i_k} \quad \text{for } 1 \leq i_1 < i_2 < \cdots < i_k \leq n$$

are a basis of $\bigwedge^k(V)$, and $\bigwedge^k(V) = 0$ when $k > n$ (when $k = 0$ the basis vector is the element $1 \in F$). In particular, $\dim_F(\bigwedge^k(V)) = \binom{n}{k}$.

Proof: As the proof of Theorem 36 shows, modulo $\mathcal{A}^k(M)$, the order of the terms in any simple k -tensor can be rearranged up to introducing a sign change. It follows that the k -tensors in the corollary (which have been arranged with increasing subscripts on the v_i and with no repeated entries) are generators for $\bigwedge^k(V)$. To show these vectors are linearly independent it suffices to exhibit an alternating k -multilinear function from V^k to F which is 1 on a given $v_{i_1} \wedge v_{i_2} \wedge \cdots \wedge v_{i_k}$ and zero on all other generators. Such a function f is defined on the basis of $\mathcal{T}^k(V)$ in Proposition 32 by $f(v_{j_1} \otimes v_{j_2} \otimes \cdots \otimes v_{j_k}) = \epsilon(\sigma)$ if σ is the unique permutation of (j_1, j_2, \dots, j_k) into (i_1, i_2, \dots, i_k) , and f is zero on every basis tensor whose k -tuple of indices cannot be permuted to (i_1, i_2, \dots, i_k) (where $\epsilon(\sigma)$ is the sign of σ). Note that f is zero on any basis tensor with repeated entries. The value $\epsilon(\sigma)$ ensures that when f is extended to all elements of $\mathcal{T}^k(V)$ it gives an alternating map, i.e., f factors through $\mathcal{A}^k(V)$. Hence f is the desired function. The computation of the dimension of $\bigwedge^k(V)$ (i.e., of the number of increasing sequences of k -tuples of indices) is left to the exercises.

The results in Corollary 37 are true for any *free* R -module of rank n . In particular if $M \cong R^n$ with R -module basis m_1, \dots, m_n then

$$\bigwedge^n(M) = R(m_1 \wedge \cdots \wedge m_n)$$

is a free (rank 1) R -module with generator $m_1 \wedge \cdots \wedge m_n$ and

$$\bigwedge^{n+1}(M) = \bigwedge^{n+2}(M) = \cdots = 0.$$

Example

Let R be the polynomial ring $\mathbb{Z}[x, y]$ in the variables x and y . If $M = R$, then $\bigwedge^2(M) = 0$ so, for example, there are no nontrivial alternating bilinear maps on $R \times R$ by the universal property of $\bigwedge^2(R)$ with respect to such maps (Theorem 36).

Suppose now that $M = I$ is the ideal (x, y) generated by x and y in R . Then $I \wedge I \neq 0$. Perhaps the easiest way to see this is to construct a nontrivial alternating bilinear map on $I \times I$. The map

$$\varphi(ax + by, cx + dy) = (ad - bc) \bmod(x, y)$$

is a well defined alternating R -bilinear map from $I \times I$ to $\mathbb{Z} = R/I$ (cf. Exercise 7). Since $\varphi(x, y) = 1$, it follows that $x \wedge y \in \bigwedge^2(I)$ is nonzero. Unlike the situation of free modules as in the examples following Theorem 36 (where arguments involving *bases* could be used), in this case it is not at all a trivial matter to give a direct verification that $x \wedge y \neq 0$ in $\bigwedge^2(I)$.

Remark: The ideal I is an example of a rank 1 (but *not free*) R -module (the rank of a module over an integral domain is defined in Section 12.1), and this example shows that the results of Corollary 37 are not true in general if the R -module is not free over R .

Homomorphisms of Tensor Algebras

If $\varphi : M \rightarrow N$ is any R -module homomorphism, then there is an induced map on the k^{th} tensor power:

$$\mathcal{T}^k(\varphi) : m_1 \otimes m_2 \otimes \cdots \otimes m_k \mapsto \varphi(m_1) \otimes \varphi(m_2) \otimes \cdots \otimes \varphi(m_k).$$

It follows directly that this map sends generators of each of the homogeneous components of the ideals $\mathcal{C}(M)$ and $\mathcal{A}(M)$ to themselves. Thus φ induces R -module homomorphisms on the quotients:

$$\mathcal{S}^k(\varphi) : \mathcal{S}^k(M) \longrightarrow \mathcal{S}^k(N) \quad \text{and} \quad \bigwedge^k(\varphi) : \bigwedge^k(M) \longrightarrow \bigwedge^k(N).$$

Moreover, each of these three maps is a ring homomorphism (hence they are graded R -algebra homomorphisms).

Of particular interest is the case when $M = V$ is an n -dimensional vector space over the field F and $\varphi : V \rightarrow V$ is an endomorphism. In this case by Corollary 37, $\bigwedge^n(\varphi)$ maps the 1-dimensional space $\bigwedge^n(V)$ to itself. Let v_1, \dots, v_n be a basis of V , so that $v_1 \wedge \cdots \wedge v_n$ is a basis of $\bigwedge^n(V)$. Then

$$\bigwedge^n(\varphi)(v_1 \wedge \cdots \wedge v_n) = \varphi(v_1) \wedge \cdots \wedge \varphi(v_n) = D(\varphi)v_1 \wedge \cdots \wedge v_n$$

for some scalar $D(\varphi) \in F$.

For any $n \times n$ matrix A over F we can define the associated endomorphism φ (with respect to the given basis v_1, \dots, v_n), which gives a map $D : M_{n \times n}(F) \rightarrow F$ where $D(A) = D(\varphi)$. It is easy to check that this map D satisfies the three axioms for a determinant function in Section 4. Then the uniqueness statement of Theorem 24 gives:

Proposition 38. If φ is an endomorphism on a n -dimensional vector space V , then $\bigwedge^n(\varphi)(w) = \det(\varphi)w$ for all $w \in \bigwedge^n(V)$.

Note that Proposition 38 characterizes the determinant of the endomorphism φ as a certain naturally induced *linear* map on $\bigwedge^n(V)$. The fact that the determinant arises naturally when considering alternating multilinear maps also explains the source of the map φ in the example above.

As with the tensor product, the maps $\mathcal{S}^k(\varphi)$ and $\bigwedge^k(\varphi)$ induced from an injective map from M to N need not remain injective (so $\bigwedge^2(M)$ need not be a submodule of $\bigwedge^2(N)$ when M is a submodule of N , for example).

Example

The inclusion $\varphi : I \hookrightarrow R$ of the ideal (x, y) into the ring $R = \mathbb{Z}[x, y]$, both considered as R -modules, induces a map

$$\bigwedge^2(\varphi) : \bigwedge^2(I) \rightarrow \bigwedge^2(R).$$

Since $\bigwedge^2(R) = 0$ and $\bigwedge^2(I) \neq 0$, the map cannot be injective.

One can show that if M is an R -module *direct summand* of N , then $\mathcal{T}(M)$ (respectively, $\mathcal{S}(M)$ and $\bigwedge(M)$) is an R -subalgebra of $\mathcal{T}(N)$ (respectively, $\mathcal{S}(N)$ and $\bigwedge(N)$) (cf. the exercises). When $R = F$ is a field then *every* subspace M of N is a direct summand of N and so the corresponding algebra for M is a subalgebra of the algebra for N .

Symmetric and Alternating Tensors

The symmetric and exterior algebras can in some instances also be defined in terms of *symmetric* and *alternating* tensors (defined below), which identify these algebras as subalgebras of the tensor algebra rather than as quotient algebras.

For any R -module M there is a natural left group action of the symmetric group S_k on $M \times M \times \cdots \times M$ (k factors) given by permuting the factors:

$$\sigma(m_1, m_2, \dots, m_k) = (m_{\sigma^{-1}(1)}, m_{\sigma^{-1}(2)}, \dots, m_{\sigma^{-1}(k)}) \quad \text{for each } \sigma \in S_k$$

(the reason for σ^{-1} is to make this a *left* group action, cf. Exercise 8 of Section 5.1). This map is clearly R -multilinear, so there is a well defined R -linear left group action of S_k on $\mathcal{T}^k(M)$ which is defined on simple tensors by

$$\sigma(m_1 \otimes m_2 \otimes \cdots \otimes m_k) = m_{\sigma^{-1}(1)} \otimes m_{\sigma^{-1}(2)} \otimes \cdots \otimes m_{\sigma^{-1}(k)} \quad \text{for each } \sigma \in S_k.$$

Definition.

- (1) An element $z \in \mathcal{T}^k(M)$ is called a *symmetric* k -tensor if $\sigma z = z$ for all σ in the symmetric group S_k .
- (2) An element $z \in \mathcal{T}^k(M)$ is called an *alternating* k -tensor if $\sigma z = \epsilon(\sigma)z$ for all σ in the symmetric group S_k , where $\epsilon(\sigma)$ is the sign, ± 1 , of the permutation σ .

It is immediate from the definition that the collection of symmetric (respectively, alternating) k -tensors is an R -submodule of the module of all k -tensors.

Example

The elements $m \otimes m$ and $m_1 \otimes m_2 + m_2 \otimes m_1$ are symmetric 2-tensors. The element $m_1 \otimes m_2 - m_2 \otimes m_1$ is an alternating 2-tensor.

It is also clear from the definition that both $\mathcal{C}^k(M)$ and $\mathcal{A}^k(M)$ are stable under the action of S_k , hence there is an induced action on the quotients $\mathcal{S}^k(M)$ and $\bigwedge^k(M)$.

Proposition 39. Let σ be an element in the symmetric group S_k and let $\epsilon(\sigma)$ be the sign of the permutation σ . Then

- (1) for every $w \in \mathcal{S}^k(M)$ we have $\sigma w = w$, and
- (2) for every $w \in \bigwedge^k(M)$ we have $\sigma w = \epsilon(\sigma)w$.

Proof: The first statement is immediate from (1) in Theorem 34. We showed in the course of the proof of Theorem 36 that

$$m_1 \wedge \cdots \wedge m_i \wedge m_{i+1} \wedge \cdots \wedge m_k = -m_1 \wedge \cdots \wedge m_{i+1} \wedge m_i \wedge \cdots \wedge m_k,$$

which shows that the formula in (2) is valid on simple products for the transposition $\sigma = (i \ i+1)$. Since these transpositions generate S_k and ϵ is a group homomorphism it follows that (2) is valid for any $\sigma \in S_k$ on simple products w . Since both sides are R -linear in w , it follows that (2) holds for all $w \in \bigwedge^k(M)$.

By Proposition 39, the symmetric group S_k acts trivially on both the submodule of symmetric k -tensors and the quotient module $\mathcal{S}^k(M)$, the k^{th} symmetric power of M . Similarly, S_k acts the same way on the submodule of alternating k -tensors as on $\bigwedge^k(M)$, the k^{th} exterior power of M . We now show that when $k!$ is a unit in R that these respective submodules and quotient modules are isomorphic (where $k!$ is the sum of the 1 of R with itself $k!$ times).

For any k -tensor $z \in \mathcal{T}^k(M)$ define

$$\begin{aligned} \text{Sym}(z) &= \sum_{\sigma \in S_k} \sigma z \\ \text{Alt}(z) &= \sum_{\sigma \in S_k} \epsilon(\sigma) \sigma z. \end{aligned}$$

For any k -tensor z , the k -tensor $\text{Sym}(z)$ is symmetric and the k -tensor $\text{Alt}(z)$ is alternating. For example, for any $\tau \in S_k$

$$\begin{aligned} \tau \text{Alt}(z) &= \sum_{\sigma \in S_k} \epsilon(\sigma) \tau \sigma z \\ &= \sum_{\sigma' \in S_k} \epsilon(\tau^{-1}\sigma') \sigma' z \quad (\text{letting } \sigma' = \tau\sigma) \\ &= \epsilon(\tau^{-1}) \sum_{\sigma' \in S_k} \epsilon(\sigma') \sigma' z = \epsilon(\tau) \text{Alt}(z). \end{aligned}$$

The tensor $\text{Sym}(z)$ is sometimes called the *symmetrization* of z and $\text{Alt}(z)$ the *skew-symmetrization* of z .

If z is already a symmetric (respectively, alternating) tensor then $\text{Sym}(z)$ (respectively, $\text{Alt}(z)$) is just $k!z$. It follows that Sym (respectively, Alt) is an R -module endomorphism of $\mathcal{T}^k(M)$ whose image lies in the submodule of symmetric (respectively, alternating) tensors. In general these maps are not surjective, but if $k!$ is a unit in R then

$$\begin{aligned} \frac{1}{k!} \text{Sym}(z) &= z \quad \text{for any symmetric tensor } z, \text{ and} \\ \frac{1}{k!} \text{Alt}(z) &= z \quad \text{for any alternating tensor } z \end{aligned}$$

so that in this case the maps $(1/k!) \text{Sym}$ and $(1/k!) \text{Alt}$ give surjective R -module homomorphisms from $\mathcal{T}^k(M)$ to the submodule of symmetric (respectively, alternating) tensors.

Proposition 40. Suppose $k!$ is a unit in the ring R and M is an R -module. Then

- (1) The map $(1/k!)Sym$ induces an R -module isomorphism between the k^{th} symmetric power of M and the R -submodule of symmetric k -tensors:

$$\frac{1}{k!} Sym : S^k(M) \cong \{\text{symmetric } k\text{-tensors}\}.$$

- (2) The map $(1/k!)Alt$ induces an R -module isomorphism between the k^{th} exterior power of M and the R -submodule of alternating k -tensors:

$$\frac{1}{k!} Alt : \bigwedge^k(M) \cong \{\text{alternating } k\text{-tensors}\}.$$

Proof: We have seen that the respective maps are surjective R -homomorphisms from $\mathcal{T}^k(M)$ so to prove the proposition it suffices to check that their kernels are $\mathcal{C}^k(M)$ and $\mathcal{A}^k(M)$, respectively. We show the first and leave the second to the exercises. It is clear that Sym is 0 on any difference of two k -tensors which differ only in the order of their factors, so $\mathcal{C}^k(M)$ is contained in the kernel of $(1/k!)Sym$ by (1) of Theorem 34. For the reverse inclusion, observe that

$$z - \frac{1}{k!} Sym(z) = \frac{1}{k!} \sum_{\sigma \in S_k} (z - \sigma z)$$

for any k -tensor z . If z is in the kernel of Sym then the left hand side of this equality is just z ; and since $z - \sigma z \in \mathcal{C}^k(M)$ for every $\sigma \in S_k$ (again by (1) of Theorem 34), it follows that $z \in \mathcal{C}^k(M)$, completing the proof.

The maps $(1/k!)Sym$ and $(1/k!)Alt$ are *projections* (cf. Exercise 11 in Section 2) onto the submodules of symmetric and antisymmetric tensors, respectively. Equivalently, if $k!$ is a unit in R , we have R -module direct sums

$$\mathcal{T}^k(M) = \ker(\pi) \oplus \text{image}(\pi)$$

for $\pi = (1/k!)Sym$ or $\pi = (1/k!)Alt$. In the former case the kernel consists of $\mathcal{C}^k(M)$ and the image is the collection of symmetric tensors (in which case $\mathcal{C}^k(M)$ is said to form an R -module *complement* to the symmetric tensors). In the latter case the kernel is $\mathcal{A}^k(M)$ and the image consists of the alternating tensors.

The R -linear left group action of S_k on $\mathcal{T}^k(M)$ makes $\mathcal{T}^k(M)$ into a module over the group ring RS_k (analogous to the formation of $F[x]$ -modules described in Section 10.1). In terms of this module structure these projections give RS_k -submodule complements to the RS_k -submodules $\mathcal{C}^k(M)$ and $\mathcal{A}^k(M)$. The “averaging” technique used to construct these maps can be used to prove a very general result (Maschke’s Theorem in Section 18.1) related to actions of finite groups on vector spaces (which is the subject of the “representation theory” of finite groups in Part VI).

If $k!$ is not invertible in R then in general we do not have such S_k -invariant direct sum decompositions so it is not in general possible to identify, for example, the k^{th} exterior power of M with the alternating k -tensors of M .

Note also that when $k!$ is invertible it is possible to *define* the k^{th} exterior power of M as the collection of alternating k -tensors (this equivalent approach is sometimes found

in the literature when the theory is developed over fields such as \mathbb{R} and \mathbb{C}). In this case the multiplication of two alternating tensors z and w is defined by first taking the product $zw = z \otimes w$ in $\mathcal{T}(M)$ and then projecting the resulting tensor into the submodule of alternating tensors. Note that the simple product of two alternating tensors need not be alternating (for example, the square of an alternating tensor is a symmetric tensor).

Example

Let V be a vector space over a field F in which $k! \neq 0$. There are many *vector space complements* to $\mathcal{A}^k(V)$ in $\mathcal{T}^k(V)$ (just extend a basis for the subspace $\mathcal{A}^k(V)$ to a basis for $\mathcal{T}^k(V)$, for example). These complements depend on choices of bases for $\mathcal{T}^k(V)$ and so are indistinguishable from each other from vector space considerations alone. The additional structure on $\mathcal{T}^k(V)$ given by the action of S_k singles out a unique complement to $\mathcal{A}^k(V)$, namely the subspace of alternating tensors in Proposition 40.

Suppose that $k! \neq 0$ in F for all $k \geq 2$ (i.e., the field F has “characteristic 0,” cf. Exercise 26 in Section 7.3), for example, $F = \mathbb{Q}$. Then the full exterior algebra $\bigwedge(V) = \bigoplus_{k \geq 0} \bigwedge^k(V)$ can be identified with the collection of tensors whose homogeneous components are alternating (with respect to the appropriate symmetric groups S_k).

Multiplication in $\bigwedge(V)$ in terms of alternating tensors is rather cumbersome, however. For example let v_1, v_2, v_3 be distinct basis vectors in V . The product of the two alternating tensors $z = v_1$ and $w = v_2 \otimes v_3 - v_3 \otimes v_2$ is obtained by first computing

$$z \otimes w = v_1 \otimes v_2 \otimes v_3 - v_1 \otimes v_3 \otimes v_2$$

in the full tensor algebra. This 3-tensor is not alternating — for example,

$$(1\ 2)(z \otimes w) = v_2 \otimes v_1 \otimes v_3 - v_3 \otimes v_1 \otimes v_2 \neq -z \otimes w$$

and also $(1\ 2\ 3)(z \otimes w) = v_3 \otimes v_1 \otimes v_2 - v_2 \otimes v_1 \otimes v_3 \neq z \otimes w$. The multiplication requires that we project this tensor into the subspace of alternating tensors. This projection is given by $(1/3!)Alt(z \otimes w)$ and an easy computation shows that

$$\begin{aligned} \frac{1}{6}Alt(z \otimes w) &= \frac{1}{3} [v_1 \otimes v_2 \otimes v_3 + v_2 \otimes v_3 \otimes v_1 + v_3 \otimes v_1 \otimes v_2 \\ &\quad - v_1 \otimes v_3 \otimes v_2 - v_2 \otimes v_1 \otimes v_3 - v_3 \otimes v_2 \otimes v_1], \end{aligned}$$

so the right hand side is the product of z and w in terms of alternating tensors. The same product in terms of the quotient algebra $\bigwedge(V)$ is simply

$$v_1 \wedge (2v_2 \wedge v_3) = 2v_1 \wedge v_2 \wedge v_3.$$

EXERCISES

In these exercises R is a commutative ring with 1 and M is an R -module; F is a field and V is a finite dimensional vector space over F .

1. Prove that if M is a cyclic R -module then $\mathcal{T}(M) = \mathcal{S}(M)$, i.e., the tensor algebra $\mathcal{T}(M)$ is commutative.
2. Fill in the details for the proof of Proposition 33 that $S/I = \bigoplus_{k=0}^{\infty} S_k/I_k$. [Show first that $S_i I_j \subseteq I_{i+j}$. Use this to show that the multiplication $(S_i/I_i)(S_j/I_j) \subseteq S_{i+j}/I_{i+j}$ is well defined, and then check the ring axioms and verify the statements made in the proof of Proposition 33.]

3. Show that the image of the map Sym_2 for the \mathbb{Z} -module \mathbb{Z} consists of the 2-tensors $a(1 \otimes 1)$ where a is an even integer. Conclude in particular that the symmetric tensor $1 \otimes 1$ in $\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Z}$ is not contained in the image of the map Sym .
4. Prove that $m \wedge n_1 \wedge n_2 \wedge \cdots \wedge n_k = (-1)^k(n_1 \wedge n_2 \wedge \cdots \wedge n_k \wedge m)$. In particular, $x \wedge (y \wedge z) = (y \wedge z) \wedge x$ for all $x, y, z \in M$.
5. Prove that if M is a free R -module of rank n then $\bigwedge^i(M)$ is a free R -module of rank $\binom{n}{i}$ for $i = 0, 1, 2, \dots$.
6. If A is any R -algebra in which $a^2 = 0$ for all $a \in A$ and $\varphi : M \rightarrow A$ is an R -module homomorphism, prove there is a unique R -algebra homomorphism $\Phi : \bigwedge(M) \rightarrow A$ such that $\Phi|_M = \varphi$.
7. Let $R = \mathbb{Z}[x, y]$ and $I = (x, y)$.
 - (a) Prove that if $ax + by = a'x + b'y$ in R then $a' = a + yf$ and $b' = b - xf$ for some polynomial $f(x, y) \in R$.
 - (b) Prove that the map $\varphi(ax + by, cx + dy) = ad - bc \bmod(x, y)$ in the example following Corollary 37 is a well defined alternating R -bilinear map from $I \times I$ to $\mathbb{Z} = R/I$.
8. Let R be an integral domain and let F be its field of fractions.
 - (a) Considering F as an R -module, prove that $\bigwedge^2 F = 0$.
 - (b) Let I be any R -submodule of F (for example, any ideal in R). Prove that $\bigwedge^i I$ is a torsion R -module for $i \geq 2$ (i.e., for every $x \in \bigwedge^i I$ there is some nonzero $r \in R$ with $rx = 0$).
 - (c) Give an example of an integral domain R and an R -module I in F with $\bigwedge^i I \neq 0$ for every $i \geq 0$ (cf. the example following Corollary 37).
9. Let $R = \mathbb{Z}[G]$ be the group ring of the group $G = \{1, \sigma\}$ of order 2. Let $M = \mathbb{Z}e_1 + \mathbb{Z}e_2$ be a free \mathbb{Z} -module of rank 2 with basis e_1 and e_2 . Define $\sigma(e_1) = e_1 + 2e_2$ and $\sigma(e_2) = -e_2$. Prove that this makes M into an R -module and that the R -module $\bigwedge^2 M$ is a group of order 2 with $e_1 \wedge e_2$ as generator.
10. Prove that $z - (1/k!)Alt(z) = (1/k!) \sum_{\sigma \in S_k} (z - \epsilon(\sigma)\sigma z)$ for any k -tensor z and use this to prove that the kernel of the R -module homomorphism $(1/k!)Alt$ in Proposition 40 is $\mathcal{A}^k(M)$.
11. Prove that the image of Alt_k is the unique largest subspace of $\mathcal{T}^k(V)$ on which each permutation σ in the symmetric group S_k acts as multiplication by the scalar $\epsilon(\sigma)$.
12. (a) Prove that if $f(x, y)$ is an alternating bilinear map on V (i.e., $f(x, x) = 0$ for all $x \in V$) then $f(x, y) = -f(y, x)$ for all $x, y \in V$.
- (b) Suppose that $-1 \neq 1$ in F . Prove that $f(x, y)$ is an alternating bilinear map on V (i.e., $f(x, x) = 0$ for all $x \in V$) if and only if $f(x, y) = -f(y, x)$ for all $x, y \in V$.
- (c) Suppose that $-1 = 1$ in F . Prove that every alternating bilinear form $f(x, y)$ on V is symmetric (i.e., $f(x, y) = f(y, x)$ for all $x, y \in V$). Prove that there is a symmetric bilinear map on V that is not alternating. [One approach: show that $C^2(V) \subset \mathcal{A}^2(V)$ and $C^2(V) \neq \mathcal{A}^2(V)$ by counting dimensions. Alternatively, construct an explicit symmetric map that is not alternating.]
13. Let F be any field in which $-1 \neq 1$ and let V be a vector space over F . Prove that $V \otimes_F V = \mathcal{S}^2(V) \oplus \bigwedge^2(V)$ i.e., that every 2-tensor may be written uniquely as a sum of a symmetric and an alternating tensor.
14. Prove that if M is an R -module *direct factor* of the R -module N then $\mathcal{T}(M)$ (respectively, $\mathcal{S}(M)$ and $\bigwedge(M)$) is an R -subalgebra of $\mathcal{T}(N)$ (respectively, $\mathcal{S}(N)$ and $\bigwedge(N)$).

CHAPTER 12

Modules over Principal Ideal Domains

The main purpose of this chapter is to prove a structure theorem for finitely generated modules over particularly nice rings, namely Principal Ideal Domains. This theorem is an example of the ideal structure of the ring (which is particularly simple for P.I.D.s) being reflected in the structure of its modules. If we apply this result in the case where the P.I.D. is the ring of integers \mathbb{Z} then we obtain a proof of the Fundamental Theorem of Finitely Generated Abelian Groups (which we examined in Chapter 5 without proof). If instead we apply this structure theorem in the case where the P.I.D. is the ring $F[x]$ of polynomials in x with coefficients in a field F we shall obtain the basic results on the so-called rational and Jordan canonical forms for a matrix. Before proceeding to the proof we briefly discuss these two important applications.

We have already discussed in Chapter 5 the result that any finitely generated abelian group is isomorphic to the direct sum of cyclic abelian groups, either \mathbb{Z} or $\mathbb{Z}/n\mathbb{Z}$ for some positive integer $n \neq 0$. Recall also that an abelian group is the same thing as a \mathbb{Z} -module. Since the ideals of \mathbb{Z} are precisely the trivial ideal (0) and the principal ideals $(n) = n\mathbb{Z}$ generated by positive integers n , we see that the Fundamental Theorem of Finitely Generated Abelian Groups in the language of modules says that any finitely generated \mathbb{Z} -module is the direct sum of modules of the form \mathbb{Z}/I where I is an ideal of \mathbb{Z} (these are the cyclic \mathbb{Z} -modules), together with a uniqueness statement when the direct sum is written in a particular form. Note the correspondence between the ideal structure of \mathbb{Z} and the structure of its (finitely generated) modules, the finitely generated abelian groups.

The Fundamental Theorem of Finitely Generated Modules over a P.I.D. states that the same result holds when the Principal Ideal Domain \mathbb{Z} is replaced by *any* P.I.D. In particular, we have seen in Chapter 10 that a module over the ring $F[x]$ of polynomials in x with coefficients in the field F is the same thing as a vector space V together with a fixed linear transformation T of V (where the element x acts on V by the linear transformation T). The Fundamental Theorem in this case will say that such a vector space is the direct sum of modules of the form $F[x]/I$ where I is an ideal of $F[x]$, hence is either the trivial ideal (0) or a principal ideal $(f(x))$ generated by some nonzero polynomial $f(x)$ (these are the cyclic $F[x]$ -modules), again with a uniqueness statement when the direct sum is written in a particular form. If this is translated back into the language of vector spaces and linear transformations we can obtain information on the

linear transformation T .

For example, suppose V is a vector space of dimension n over F and we choose a basis for V . Then giving a linear transformation T of V to itself is the same thing as giving an $n \times n$ matrix A with coefficients in F (and choosing a different basis for V gives a different matrix B for T which is similar to A i.e., is of the form $P^{-1}AP$ for some invertible matrix P which defines the change of basis). We shall see that the Fundamental Theorem in this situation implies (under the assumption that the field F contains all the “eigenvalues” for the given linear transformation T) that there is a basis for V so that the associated matrix for T is *as close to being a diagonal matrix as possible* and so has a particularly simple form. This is the *Jordan canonical form*. The *rational canonical form* is another simple form for the matrix for T (that does not require the eigenvalues for T to be elements of F). In this way we shall be able to give canonical forms for arbitrary $n \times n$ matrices over fields F , that is, find matrices which are similar to a given $n \times n$ matrix and which are particularly simple (almost diagonal, for example).

Example

Let $V = \mathbb{Q}^3 = \{(x, y, z) \mid x, y, z \in \mathbb{Q}\}$ be the usual 3-dimensional vector space of ordered 3-tuples with entries from the field $F = \mathbb{Q}$ of rational numbers and suppose T is the linear transformation

$$T(x, y, z) = (9x + 4y + 5z, -4x - 3z, -6x - 4y - 2z), \quad x, y, z \in \mathbb{Q}.$$

If we take the standard basis $e_1 = (1, 0, 0)$, $e_2 = (0, 1, 0)$, $e_3 = (0, 0, 1)$ for V then the matrix A representing this linear transformation is

$$A = \begin{pmatrix} 9 & 4 & 5 \\ -4 & 0 & -3 \\ -6 & -4 & -2 \end{pmatrix}.$$

We shall see that the Jordan canonical form for this matrix A is the much simpler matrix

$$B = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}$$

obtained by taking instead the basis $f_1 = (2, -1, -2)$, $f_2 = (1, 0, -1)$, $f_3 = (3, -2, -2)$ for V , since in this case

$$\begin{aligned} T(f_1) &= T(2, -1, -2) = (4, -2, -4) = 2 \cdot f_1 + 0 \cdot f_2 + 0 \cdot f_3 \\ T(f_2) &= T(1, 0, -1) = (4, -1, -4) = 1 \cdot f_1 + 2 \cdot f_2 + 0 \cdot f_3 \\ T(f_3) &= T(3, -2, -2) = (9, -6, -6) = 0 \cdot f_1 + 0 \cdot f_2 + 3 \cdot f_3, \end{aligned}$$

so the columns of the matrix representing T with respect to this basis are $(2, 0, 0)$, $(1, 2, 0)$ and $(0, 0, 3)$, i.e., T has matrix B with respect to this basis. In particular A is similar to the simpler matrix B .

In fact this linear transformation T *cannot* be diagonalized (i.e., there is no choice of basis for V for which the corresponding matrix is a diagonal matrix) so that the matrix B is as close to a diagonal matrix for T as is possible.

The first section below gives some general definitions and states and proves the Fundamental Theorem over an arbitrary P.I.D., after which we return to the application to canonical forms (the application to abelian groups appears in Chapter 5). These applications can be read independently of the general proof. An alternate and computationally useful proof valid for Euclidean Domains (so in particular for the rings \mathbb{Z} and $F[x]$) along the lines of row and column operations is outlined in the exercises.

12.1 THE BASIC THEORY

We first describe some general finiteness conditions. Let R be a ring and let M be a left R -module.

Definition.

- (1) The left R -module M is said to be a *Noetherian R -module* or to satisfy the *ascending chain condition on submodules* (or *A.C.C. on submodules*) if there are no infinite increasing chains of submodules, i.e., whenever

$$M_1 \subseteq M_2 \subseteq M_3 \subseteq \dots$$

is an increasing chain of submodules of M , then there is a positive integer m such that for all $k \geq m$, $M_k = M_m$ (so the chain becomes stationary at stage m : $M_m = M_{m+1} = M_{m+2} = \dots$).

- (2) The ring R is said to be *Noetherian* if it is Noetherian as a left module over itself, i.e., if there are no infinite increasing chains of left ideals in R .

One can formulate analogous notions of A.C.C. on right and on two-sided ideals in a (possibly noncommutative) ring R . For noncommutative rings these properties need not be related.

Theorem 1. Let R be a ring and let M be a left R -module. Then the following are equivalent:

- (1) M is a Noetherian R -module.
- (2) Every nonempty set of submodules of M contains a maximal element under inclusion.
- (3) Every submodule of M is finitely generated.

Proof: [(1) implies (2)] Assume M is Noetherian and let Σ be any nonempty collection of submodules of M . Choose any $M_1 \in \Sigma$. If M_1 is a maximal element of Σ , (2) holds, so assume M_1 is not maximal. Then there is some $M_2 \in \Sigma$ such that $M_1 \subset M_2$. If M_2 is maximal in Σ , (2) holds, so we may assume there is an $M_3 \in \Sigma$ properly containing M_2 . Proceeding in this way one sees that if (2) fails we can produce by the Axiom of Choice an infinite strictly increasing chain of elements of Σ , contrary to (1).

[(2) implies (3)] Assume (2) holds and let N be any submodule of M . Let Σ be the collection of all finitely generated submodules of N . Since $\{0\} \in \Sigma$, this collection is nonempty. By (2) Σ contains a maximal element N' . If $N' \neq N$, let $x \in N - N'$. Since $N' \in \Sigma$, the submodule N' is finitely generated by assumption, hence also the

submodule generated by N' and x is finitely generated. This contradicts the maximality of N' , so $N = N'$ is finitely generated.

[(3) implies (1)] Assume (3) holds and let $M_1 \subseteq M_2 \subseteq M_3 \dots$ be a chain of submodules of M . Let

$$N = \bigcup_{i=1}^{\infty} M_i$$

and note that N is a submodule. By (3) N is finitely generated by, say, a_1, a_2, \dots, a_n . Since $a_i \in N$ for all i , each a_i lies in one of the submodules in the chain, say M_{j_i} . Let $m = \max\{j_1, j_2, \dots, j_n\}$. Then $a_i \in M_m$ for all i so the module they generate is contained in M_m , i.e., $N \subseteq M_m$. This implies $M_m = N = M_k$ for all $k \geq m$, which proves (1).

Corollary 2. If R is a P.I.D. then every nonempty set of ideals of R has a maximal element and R is a Noetherian ring.

Proof: The P.I.D. R satisfies condition (3) in the theorem with $M = R$.

Recall that even if M itself is a finitely generated R -module, submodules of M need not be finitely generated, so the condition that M be a Noetherian R -module is in general stronger than the condition that M be a finitely generated R -module.

We require a result on “linear dependence” before turning to the main results of this chapter.

Proposition 3. Let R be an integral domain and let M be a free R -module of rank $n < \infty$. Then any $n + 1$ elements of M are R -linearly dependent, i.e., for any $y_1, y_2, \dots, y_{n+1} \in M$ there are elements $r_1, r_2, \dots, r_{n+1} \in R$, not all zero, such that

$$r_1y_1 + r_2y_2 + \dots + r_{n+1}y_{n+1} = 0.$$

Proof: The quickest way of proving this is to embed R in its quotient field F (since R is an integral domain) and observe that since $M \cong R \oplus R \oplus \dots \oplus R$ (n times) we obtain $M \subseteq F \oplus F \oplus \dots \oplus F$. The latter is an n -dimensional vector space over F so any $n + 1$ elements of M are F -linearly dependent. By clearing the denominators of the scalars (by multiplying through by the product of all the denominators, for example), we obtain an R -linear dependence relation among the $n + 1$ elements of M .

Alternatively, let e_1, \dots, e_n be a basis of the free R -module M and let y_1, \dots, y_{n+1} be any $n + 1$ elements of M . For $1 \leq i \leq n + 1$ write $y_i = a_{1i}e_1 + a_{2i}e_2 + \dots + a_{ni}e_n$ in terms of the basis e_1, e_2, \dots, e_n . Let A be the $(n + 1) \times (n + 1)$ matrix whose i, j entry is a_{ij} , $1 \leq i \leq n$, $1 \leq j \leq n + 1$ and whose last row is zero, so certainly $\det A = 0$. Since R is an integral domain, Corollary 27 of Section 11.4 shows that the columns of A are R -linearly dependent. Any dependence relation on the columns of A gives a dependence relation on the y_i 's, completing the proof.

If R is any integral domain and M is any R -module recall that

$$\text{Tor}(M) = \{x \in M \mid rx = 0 \text{ for some nonzero } r \in R\}$$

is a submodule of M (called *the* torsion submodule of M) and if N is any submodule of $\text{Tor}(M)$, N is called *a* torsion submodule of M (so the torsion submodule of M is the union of all torsion submodules of M , i.e., is the maximal torsion submodule of M). If $\text{Tor}(M) = 0$, the module M is said to be *torsion free*.

For any submodule N of M , the *annihilator* of N is the ideal of R defined by

$$\text{Ann}(N) = \{r \in R \mid rn = 0 \text{ for all } n \in N\}.$$

Note that if N is not a torsion submodule of M then $\text{Ann}(N) = (0)$. It is easy to see that if N, L are submodules of M with $N \subseteq L$, then $\text{Ann}(L) \subseteq \text{Ann}(N)$. If R is a P.I.D. and $N \subseteq L \subseteq M$ with $\text{Ann}(N) = (a)$ and $\text{Ann}(L) = (b)$, then $a \mid b$. In particular, the annihilator of any element x of M divides the annihilator of M (this is implied by Lagrange's Theorem when $R = \mathbb{Z}$).

Definition. For any integral domain R the *rank* of an R -module M is the maximum number of R -linearly independent elements of M .

The preceding proposition states that for a free R -module M over an integral domain the rank of a submodule is bounded by the rank of M . This notion of rank agrees with previous uses of the same term. If the ring $R = F$ is a field, then the rank of an R -module M is the dimension of M as a vector space over F and any maximal set of F -linearly independent elements is a basis for M . For a general integral domain, however, an R -module M of rank n need not have a “basis,” i.e., need not be a *free* R -module even if M is torsion free, so some care is necessary with the notion of rank, particularly with respect to the torsion elements of M . Exercises 1 to 6 and 20 give an alternate characterization of the rank and provide some examples of (torsion free) R -modules (of rank 1) that are not free.

The next important result shows that if N is a submodule of a free module of finite rank over a P.I.D. then N is again a free module of finite rank and furthermore it is possible to choose generators for the two modules which are related in a simple way.

Theorem 4. Let R be a Principal Ideal Domain, let M be a free R -module of finite rank n and let N be a submodule of M . Then

- (1) N is free of rank m , $m \leq n$ and
- (2) there exists a basis y_1, y_2, \dots, y_n of M so that $a_1 y_1, a_2 y_2, \dots, a_m y_m$ is a basis of N where a_1, a_2, \dots, a_m are nonzero elements of R with the divisibility relations

$$a_1 \mid a_2 \mid \cdots \mid a_m.$$

Proof: The theorem is trivial for $N = \{0\}$, so assume $N \neq \{0\}$. For each R -module homomorphism φ of M into R , the image $\varphi(N)$ of N is a submodule of R , i.e., an ideal in R . Since R is a P.I.D. this ideal must be principal, say $\varphi(N) = (a_\varphi)$, for some $a_\varphi \in R$. Let

$$\Sigma = \{(a_\varphi) \mid \varphi \in \text{Hom}_R(M, R)\}$$

be the collection of the principal ideals in R obtained in this way from the R -module homomorphisms of M into R . The collection Σ is certainly nonempty since taking φ

to be the trivial homomorphism shows that $(0) \in \Sigma$. By Corollary 2, Σ has at least one maximal element i.e., there is at least one homomorphism v of M to R so that the principal ideal $v(N) = (a_1)$ is not properly contained in any other element of Σ . Let $a_1 = a_v$ for this maximal element and let $y \in N$ be an element mapping to the generator a_1 under the homomorphism v : $v(y) = a_1$.

We now show the element a_1 is nonzero. Let x_1, x_2, \dots, x_n be any basis of the free module M and let $\pi_i \in \text{Hom}_R(M, R)$ be the natural projection homomorphism onto the i^{th} coordinate with respect to this basis. Since $N \neq \{0\}$, there exists an i such that $\pi_i(N) \neq 0$, which in particular shows that Σ contains more than just the trivial ideal (0) . Since (a_1) is a maximal element of Σ it follows that $a_1 \neq 0$.

We next show that this element a_1 divides $\varphi(y)$ for every $\varphi \in \text{Hom}_R(M, R)$. To see this let d be a generator for the principal ideal generated by a_1 and $\varphi(y)$. Then d is a divisor of both a_1 and $\varphi(y)$ in R and $d = r_1 a_1 + r_2 \varphi(y)$ for some $r_1, r_2 \in R$. Consider the homomorphism $\psi = r_1 v + r_2 \varphi$ from M to R . Then $\psi(y) = (r_1 v + r_2 \varphi)(y) = r_1 a_1 + r_2 \varphi(y) = d$ so that $d \in \psi(N)$, hence also $(d) \subseteq \psi(N)$. But d is a divisor of a_1 so we also have $(a_1) \subseteq (d)$. Then $(a_1) \subseteq (d) \subseteq \psi(N)$ and by the maximality of (a_1) we must have equality: $(a_1) = (d) = \psi(N)$. In particular $(a_1) = (d)$ shows that $a_1 \mid \varphi(y)$ since d divides $\varphi(y)$.

If we apply this to the projection homomorphisms π_i we see that a_1 divides $\pi_i(y)$ for all i . Write $\pi_i(y) = a_1 b_i$ for some $b_i \in R$, $1 \leq i \leq n$ and define

$$y_1 = \sum_{i=1}^n b_i x_i.$$

Note that $a_1 y_1 = y$. Since $a_1 = v(y) = v(a_1 y_1) = a_1 v(y_1)$ and a_1 is a nonzero element of the integral domain R this shows

$$v(y_1) = 1.$$

We now verify that this element y_1 can be taken as one element in a basis for M and that $a_1 y_1$ can be taken as one element in a basis for N , namely that we have

(a) $M = Ry_1 \oplus \ker v$, and

(b) $N = Ra_1 y_1 \oplus (N \cap \ker v)$.

To see (a) let x be an arbitrary element in M and write $x = v(x)y_1 + (x - v(x)y_1)$. Since

$$\begin{aligned} v(x - v(x)y_1) &= v(x) - v(x)v(y_1) \\ &= v(x) - v(x) \cdot 1 \\ &= 0 \end{aligned}$$

we see that $x - v(x)y_1$ is an element in the kernel of v . This shows that x can be written as the sum of an element in Ry_1 and an element in the kernel of v , so $M = Ry_1 + \ker v$. To see that the sum is direct, suppose ry_1 is also an element in the kernel of v . Then $0 = v(ry_1) = rv(y_1) = r$ shows that this element is indeed 0.

For (b) observe that $v(x')$ is divisible by a_1 for every $x' \in N$ by the definition of a_1 as a generator for $v(N)$. If we write $v(x') = ba_1$ where $b \in R$ then the decomposition we used in (a) above is $x' = v(x')y_1 + (x' - v(x')y_1) = ba_1 y_1 + (x' - ba_1 y_1)$ where the second summand is in the kernel of v and is an element of N . This shows that

$N = Ra_1y_1 + (N \cap \ker \nu)$. The fact that the sum in (b) is direct is a special case of the directness of the sum in (a).

We now prove part (1) of the theorem by induction on the rank, m , of N . If $m = 0$, then N is a torsion module, hence $N = 0$ since a free module is torsion free, so (1) holds trivially. Assume then that $m > 0$. Since the sum in (b) above is direct we see easily that $N \cap \ker \nu$ has rank $m - 1$ (cf. Exercise 3). By induction $N \cap \ker \nu$ is then a free R -module of rank $m - 1$. Again by the directness of the sum in (b) we see that adjoining a_1y_1 to any basis of $N \cap \ker \nu$ gives a basis of N , so N is also free (of rank m), which proves (1).

Finally, we prove (2) by induction on n , the rank of M . Applying (1) to the submodule $\ker \nu$ shows that this submodule is free and because the sum in (a) is direct it is free of rank $n - 1$. By the induction assumption applied to the module $\ker \nu$ (which plays the role of M) and its submodule $\ker \nu \cap N$ (which plays the role of N), we see that there is a basis y_2, y_3, \dots, y_n of $\ker \nu$ such that $a_2y_2, a_3y_3, \dots, a_my_m$ is a basis of $N \cap \ker \nu$ for some elements a_2, a_3, \dots, a_m of R with $a_2 \mid a_3 \mid \dots \mid a_m$. Since the sums (a) and (b) are direct, y_1, y_2, \dots, y_n is a basis of M and $a_1y_1, a_2y_2, \dots, a_my_m$ is a basis of N . To complete the induction it remains to show that a_1 divides a_2 . Define a homomorphism φ from M to R by defining $\varphi(y_1) = \varphi(y_2) = 1$ and $\varphi(y_i) = 0$, for all $i > 2$, on the basis for M . Then for this homomorphism φ we have $a_1 = \varphi(a_1y_1)$ so $a_1 \in \varphi(N)$ hence also $(a_1) \subseteq \varphi(N)$. By the maximality of (a_1) in Σ it follows that $(a_1) = \varphi(N)$. Since $a_2 = \varphi(a_2y_2) \in \varphi(N)$ we then have $a_2 \in (a_1)$ i.e., $a_1 \mid a_2$. This completes the proof of the theorem.

Recall that the left R -module C is a *cyclic* R -module (for any ring R , not necessarily commutative nor with 1) if there is an element $x \in C$ such that $C = Rx$. We can then define an R -module homomorphism

$$\pi : R \rightarrow C$$

by $\pi(r) = rx$, which will be surjective by the assumption $C = Rx$. The First Isomorphism Theorem gives an isomorphism of (left) R -modules

$$R/\ker \pi \cong C.$$

If R is a P.I.D., $\ker \pi$ is a principal ideal, (a) , so we see that the cyclic R -modules C are of the form $R/(a)$ where $(a) = \text{Ann}(C)$.

The cyclic modules are the simplest modules (since they require only one generator). The existence portion of the Fundamental Theorem states that any finitely generated module over a P.I.D. is isomorphic to the direct sum of finitely many cyclic modules.

Theorem 5. (Fundamental Theorem, Existence: Invariant Factor Form) Let R be a P.I.D. and let M be a finitely generated R -module.

- (1) Then M is isomorphic to the direct sum of finitely many cyclic modules. More precisely,

$$M \cong R^r \oplus R/(a_1) \oplus R/(a_2) \oplus \cdots \oplus R/(a_m)$$

for some integer $r \geq 0$ and nonzero elements a_1, a_2, \dots, a_m of R which are not units in R and which satisfy the divisibility relations

$$a_1 \mid a_2 \mid \cdots \mid a_m.$$

- (2) M is torsion free if and only if M is free.
(3) In the decomposition in (1),

$$\text{Tor}(M) \cong R/(a_1) \oplus R/(a_2) \oplus \cdots \oplus R/(a_m).$$

In particular M is a torsion module if and only if $r = 0$ and in this case the annihilator of M is the ideal (a_m) .

Proof: The module M can be generated by a finite set of elements by assumption so let x_1, x_2, \dots, x_n be a set of generators of M of minimal cardinality. Let R^n be the free R -module of rank n with basis b_1, b_2, \dots, b_n and define the homomorphism $\pi : R^n \rightarrow M$ by defining $\pi(b_i) = x_i$ for all i , which is automatically surjective since x_1, \dots, x_n generate M . By the First Isomorphism Theorem for modules we have $R^n / \ker \pi \cong M$. Now, by Theorem 4 applied to R^n and the submodule $\ker \pi$ we can choose another basis y_1, y_2, \dots, y_n of R^n so that $a_1 y_1, a_2 y_2, \dots, a_m y_m$ is a basis of $\ker \pi$ for some elements a_1, a_2, \dots, a_m of R with $a_1 \mid a_2 \mid \cdots \mid a_m$. This implies

$$M \cong R^n / \ker \pi = (Ry_1 \oplus Ry_2 \oplus \cdots \oplus Ry_n) / (Ra_1 y_1 \oplus Ra_2 y_2 \oplus \cdots \oplus Ra_m y_m).$$

To identify the quotient on the right hand side we use the natural surjective R -module homomorphism

$$Ry_1 \oplus Ry_2 \oplus \cdots \oplus Ry_n \rightarrow R/(a_1) \oplus R/(a_2) \oplus \cdots \oplus R/(a_m) \oplus R^{n-m}$$

that maps $(\alpha_1 y_1, \dots, \alpha_n y_n)$ to $(\alpha_1 \bmod (a_1), \dots, \alpha_m \bmod (a_m), \alpha_{m+1}, \dots, \alpha_n)$. The kernel of this map is clearly the set of elements where a_i divides α_i , $i = 1, 2, \dots, m$, i.e., $Ra_1 y_1 \oplus Ra_2 y_2 \oplus \cdots \oplus Ra_m y_m$ (cf. Exercise 7). Hence we obtain

$$M \cong R/(a_1) \oplus R/(a_2) \oplus \cdots \oplus R/(a_m) \oplus R^{n-m}.$$

If a is a unit in R then $R/(a) = 0$, so in this direct sum we may remove any of the initial a_i which are units. This gives the decomposition in (1) (with $r = n - m$).

Since $R/(a)$ is a torsion R -module for any nonzero element a of R , (1) immediately implies M is a torsion free module if and only if $M \cong R^r$, which is (2). Part (3) is immediate from the definitions since the annihilator of $R/(a)$ is evidently the ideal (a) .

We shall shortly prove the uniqueness of the decomposition in Theorem 5, namely that if we have

$$M \cong R^{r'} \oplus R/(b_1) \oplus R/(b_2) \oplus \cdots \oplus R/(b_{m'})$$

for some integer $r' \geq 0$ and nonzero elements $b_1, b_2, \dots, b_{m'}$ of R which are not units with

$$b_1 \mid b_2 \mid \cdots \mid b_{m'},$$

then $r = r'$, $m = m'$ and $(a_i) = (b_i)$ (so $a_i = b_i$ up to units) for all i . It is precisely the divisibility condition $a_1 \mid a_2 \mid \cdots \mid a_m$ which gives this uniqueness.

Definition. The integer r in Theorem 5 is called the *free rank* or the *Betti number* of M and the elements $a_1, a_2, \dots, a_m \in R$ (defined up to multiplication by units in R) are called the *invariant factors* of M .

Note that until we have proved that the invariant factors of M are unique we should properly refer to a set of invariant factors for M (and similarly for the free rank), by which we mean any elements giving a decomposition for M as in (1) of the theorem above.

Using the Chinese Remainder Theorem it is possible to decompose the cyclic modules in Theorem 5 further so that M is the direct sum of cyclic modules whose annihilators are as simple as possible (namely (0) or generated by powers of primes in R). This gives an alternate decomposition which we shall also see is unique and which we now describe.

Suppose a is a nonzero element of the Principal Ideal Domain R . Then since R is also a Unique Factorization Domain we can write

$$a = up_1^{\alpha_1} p_2^{\alpha_2} \cdots p_s^{\alpha_s}$$

where the p_i are distinct primes in R and u is a unit. This factorization is unique up to units, so the ideals $(p_i^{\alpha_i})$, $i = 1, \dots, s$ are uniquely defined. For $i \neq j$ we have $(p_i^{\alpha_i}) + (p_j^{\alpha_j}) = R$ since the sum of these two ideals is generated by a greatest common divisor, which is 1 for distinct primes p_i, p_j . Put another way, the ideals $(p_i^{\alpha_i})$, $i = 1, \dots, s$, are comaximal in pairs. The intersection of all these ideals is the ideal (a) since a is the least common multiple of $p_1^{\alpha_1}, p_2^{\alpha_2}, \dots, p_s^{\alpha_s}$. Then the Chinese Remainder Theorem (Theorem 7.17) shows that

$$R/(a) \cong R/(p_1^{\alpha_1}) \oplus R/(p_2^{\alpha_2}) \oplus \cdots \oplus R/(p_s^{\alpha_s})$$

as rings and also as R -modules.

Applying this to the modules in Theorem 5 allows us to write each of the direct summands $R/(a_i)$ for the invariant factor a_i of M as a direct sum of cyclic modules whose annihilators are the prime power divisors of a_i . This proves:

Theorem 6. (Fundamental Theorem, Existence: Elementary Divisor Form) Let R be a P.I.D. and let M be a finitely generated R -module. Then M is the direct sum of a finite number of cyclic modules whose annihilators are either (0) or generated by powers of primes in R , i.e.,

$$M \cong R^r \oplus R/(p_1^{\alpha_1}) \oplus R/(p_2^{\alpha_2}) \oplus \cdots \oplus R/(p_t^{\alpha_t})$$

where $r \geq 0$ is an integer and $p_1^{\alpha_1}, \dots, p_t^{\alpha_t}$ are positive powers of (not necessarily distinct) primes in R .

We proved Theorem 6 by using the prime power factors of the invariant factors for M . In fact we shall see that the decomposition of M into a direct sum of cyclic modules whose annihilators are (0) or prime powers as in Theorem 6 is unique, i.e., the integer r and the ideals $(p_1^{\alpha_1}), \dots, (p_t^{\alpha_t})$ are uniquely defined for M . These prime powers are given a name:

Definition. Let R be a P.I.D. and let M be a finitely generated R -module as in Theorem 6. The prime powers $p_1^{\alpha_1}, \dots, p_t^{\alpha_t}$ (defined up to multiplication by units in R) are called the *elementary divisors* of M .

Suppose M is a finitely generated torsion module over the Principal Ideal Domain R . If for the *distinct* primes p_1, p_2, \dots, p_n occurring in the decomposition in Theorem 6 we group together all the cyclic factors corresponding to the same prime p_i we see in particular that M can be written as a direct sum

$$M = N_1 \oplus N_2 \oplus \cdots \oplus N_n$$

where N_i consists of all the elements of M which are annihilated by some power of the prime p_i . This result holds also for modules over R which may not be finitely generated:

Theorem 7. (The Primary Decomposition Theorem) Let R be a P.I.D. and let M be a nonzero torsion R -module (not necessarily finitely generated) with nonzero annihilator a . Suppose the factorization of a into distinct prime powers in R is

$$a = up_1^{\alpha_1} p_2^{\alpha_2} \cdots p_n^{\alpha_n}$$

and let $N_i = \{x \in M \mid p_i^{\alpha_i}x = 0\}$, $1 \leq i \leq n$. Then N_i is a submodule of M with annihilator $p_i^{\alpha_i}$ and is the submodule of M of all elements annihilated by some power of p_i . We have

$$M = N_1 \oplus N_2 \oplus \cdots \oplus N_n.$$

If M is finitely generated then each N_i is the direct sum of finitely many cyclic modules whose annihilators are divisors of $p_i^{\alpha_i}$.

Proof: We have already proved these results in the case where M is finitely generated over R . In the general case it is clear that N_i is a submodule of M with annihilator dividing $p_i^{\alpha_i}$. Since R is a P.I.D. the ideals $(p_i^{\alpha_i})$ and $(p_j^{\alpha_j})$ are comaximal for $i \neq j$, so the direct sum decomposition of M can be proved easily by modifying the argument in the proof of the Chinese Remainder Theorem to apply it to modules. Using this direct sum decomposition it is easy to see that the annihilator of N_i is precisely $p_i^{\alpha_i}$.

Definition. The submodule N_i in the previous theorem is called the p_i -primary component of M .

Notice that with this terminology the elementary divisors of a finitely generated module M are just the invariant factors of the primary components of $\text{Tor}(M)$.

We now prove the uniqueness statements regarding the decompositions in the Fundamental Theorem.

Note that if M is any module over a commutative ring R and a is an element of R then $aM = \{am \mid m \in M\}$ is a submodule of M . Recall also that in a Principal Ideal Domain R the nonzero prime ideals are maximal, hence the quotient of R by a nonzero prime ideal is a field.

Lemma 8. Let R be a P.I.D. and let p be a prime in R . Let F denote the field $R/(p)$.

- (1) Let $M = R^r$. Then $M/pM \cong F^r$.
- (2) Let $M = R/(a)$ where a is a nonzero element of R . Then

$$M/pM \cong \begin{cases} F & \text{if } p \text{ divides } a \text{ in } R \\ 0 & \text{if } p \text{ does not divide } a \text{ in } R. \end{cases}$$

- (3) Let $M = R/(a_1) \oplus R/(a_2) \oplus \cdots \oplus R/(a_k)$ where each a_i is divisible by p . Then $M/pM \cong F^k$.

Proof: (1) There is a natural map from R^r to $(R/(p))^r$ defined by mapping $(\alpha_1, \dots, \alpha_r)$ to $(\alpha_1 \bmod (p), \dots, \alpha_r \bmod (p))$. This is clearly a surjective R -module homomorphism with kernel consisting of the r -tuples all of whose coordinates are divisible by p , i.e., pR^r , so $R^r/pR^r \cong (R/(p))^r$, which is (1).

(2) This follows from the Isomorphism Theorems: note first that $p(R/(a))$ is the image of the ideal (p) in the quotient $R/(a)$, hence is $(p)+(a)/(a)$. The ideal $(p)+(a)$ is generated by a greatest common divisor of p and a , hence is (p) if p divides a and is $R = (1)$ otherwise. Hence $pM = (p)/(a)$ if p divides a and is $R/(a) = M$ otherwise. If p divides a then $M/pM = (R/(a))/((p)/(a)) \cong R/(p)$, and if p does not divide a then $M/pM = M/M = 0$, which proves (2).

(3) This follows from (2) as in the proof of part (1) of Theorem 5.

Theorem 9. (Fundamental Theorem, Uniqueness) Let R be a P.I.D.

- (1) Two finitely generated R -modules M_1 and M_2 are isomorphic if and only if they have the same free rank and the same list of invariant factors.
- (2) Two finitely generated R -modules M_1 and M_2 are isomorphic if and only if they have the same free rank and the same list of elementary divisors.

Proof: If M_1 and M_2 have the same free rank and list of invariant factors or the same free rank and list of elementary divisors then they are clearly isomorphic.

Suppose that M_1 and M_2 are isomorphic. Any isomorphism between M_1 and M_2 maps the torsion in M_1 to the torsion in M_2 so we must have $\text{Tor}(M_1) \cong \text{Tor}(M_2)$. Then $R^{r_1} \cong M_1/\text{Tor}(M_1) \cong M_2/\text{Tor}(M_2) \cong R^{r_2}$ where r_1 is the free rank of M_1 and r_2 is the free rank of M_2 . Let p be any nonzero prime in R . Then from $R^{r_1} \cong R^{r_2}$ we obtain $R^{r_1}/pR^{r_1} \cong R^{r_2}/pR^{r_2}$. By (1) of the previous lemma, this implies $F^{r_1} \cong F^{r_2}$ where F is the field R/pR . Hence we have an isomorphism of an r_1 -dimensional vector space over F with an r_2 -dimensional vector space over F , so that $r_1 = r_2$ and M_1 and M_2 have the same free rank.

We are reduced to showing that M_1 and M_2 have the same lists of invariant factors and elementary divisors. To do this we need only work with the isomorphic torsion modules $\text{Tor}(M_1)$ and $\text{Tor}(M_2)$, i.e., we may as well assume that both M_1 and M_2 are torsion R -modules.

We first show they have the same elementary divisors. It suffices to show that for any fixed prime p the elementary divisors which are a power of p are the same for both M_1 and M_2 . If $M_1 \cong M_2$ then the p -primary submodule of M_1 (= the direct

sum of the cyclic factors whose elementary divisors are powers of p) is isomorphic to the p -primary submodule of M_2 , since these are the submodules of elements which are annihilated by some power of p . We are therefore reduced to the case of proving that if two modules M_1 and M_2 which have annihilator a power of p are isomorphic then they have the same elementary divisors.

We proceed by induction on the power of p in the annihilator of M_1 (which is the same as the annihilator of M_2 since M_1 and M_2 are isomorphic). If this power is 0, then both M_1 and M_2 are 0 and we are done. Otherwise M_1 (and M_2) have nontrivial elementary divisors. Suppose the elementary divisors of M_1 are given by

$$\text{elementary divisors of } M_1: \underbrace{p, p, \dots, p}_{m \text{ times}}, p^{\alpha_1}, p^{\alpha_2}, \dots, p^{\alpha_s},$$

where $2 \leq \alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_s$, i.e., M_1 is the direct sum of cyclic modules with generators $x_1, x_2, \dots, x_m, x_{m+1}, \dots, x_{m+s}$, say, whose annihilators are $(p), (p), \dots, (p), (p^{\alpha_1}), \dots, (p^{\alpha_s})$, respectively. Then the submodule pM_1 has elementary divisors

$$\text{elementary divisors of } pM_1: p^{\alpha_1-1}, p^{\alpha_2-1}, \dots, p^{\alpha_s-1}$$

since pM_1 is the direct sum of the cyclic modules with generators $px_1, px_2, \dots, px_m, px_{m+1}, \dots, px_{m+s}$ whose annihilators are $(1), (1), \dots, (1), (p^{\alpha_1-1}), \dots, (p^{\alpha_s-1})$, respectively. Similarly, if the elementary divisors of M_2 are given by

$$\text{elementary divisors of } M_2: \underbrace{p, p, \dots, p}_{n \text{ times}}, p^{\beta_1}, p^{\beta_2}, \dots, p^{\beta_t},$$

where $2 \leq \beta_1 \leq \beta_2 \leq \dots \leq \beta_t$, then pM_2 has elementary divisors

$$\text{elementary divisors of } pM_2: p^{\beta_1-1}, p^{\beta_2-1}, \dots, p^{\beta_t-1}.$$

Since $M_1 \cong M_2$, also $pM_1 \cong pM_2$ and the power of p in the annihilator of pM_1 is one less than the power of p in the annihilator of M_1 . By induction, the elementary divisors for pM_1 are the same as the elementary divisors for pM_2 , i.e., $s = t$ and $\alpha_i - 1 = \beta_i - 1$ for $i = 1, 2, \dots, s$, hence $\alpha_i = \beta_i$ for $i = 1, 2, \dots, s$. Finally, since also $M_1/pM_1 \cong M_2/pM_2$ we see from (3) of the lemma above that $F^{m+s} \cong F^{n+t}$, which shows that $m + s = n + t$ hence $m = n$ since we have already seen $s = t$. This proves that the set of elementary divisors for M_1 is the same as the set of elementary divisors for M_2 .

We now show that M_1 and M_2 must have the same invariant factors. Suppose $a_1 | a_2 | \dots | a_m$ are invariant factors for M_1 . We obtain a set of elementary divisors for M_1 by taking the prime power factors of these elements. Note that then the divisibility relations on the invariant factors imply that a_m is the product of the largest of the prime powers among these elementary divisors, a_{m-1} is the product of the largest prime powers among these elementary divisors once the factors for a_m have been removed, and so on. If $b_1 | b_2 | \dots | b_n$ are invariant factors for M_2 then we similarly obtain a set of elementary divisors for M_2 by taking the prime power factors of these elements. But we showed above that the elementary divisors for M_1 and M_2 are the same, and it follows that the same is true of the invariant factors.

Corollary 10. Let R be a P.I.D. and let M be a finitely generated R -module.

- (1) The elementary divisors of M are the prime power factors of the invariant factors of M .
- (2) The largest invariant factor of M is the product of the largest of the distinct prime powers among the elementary divisors of M , the next largest invariant factor is the product of the largest of the distinct prime powers among the remaining elementary divisors of M , and so on.

Proof: The procedure in (1) gives a set of elementary divisors and since the elementary divisors for M are unique by the theorem, it follows that the procedure in (1) gives *the* set of elementary divisors. Similarly for (2).

Corollary 11. (The Fundamental Theorem of Finitely Generated Abelian Groups) See Theorem 5.3 and Theorem 5.5.

Proof: Take $R = \mathbb{Z}$ in Theorems 5, 6 and 9 (note however that the invariant factors are listed in reverse order in Chapter 5 for computational convenience).

The procedure for passing between elementary divisors and invariant factors in Corollary 10 is described in some detail in Chapter 5 in the case of finitely generated abelian groups.

Note also that if a finitely generated module M is written as a direct sum of cyclic modules of the form $R/(a)$ then the ideals (a) which occur are not in general unique unless some additional conditions are imposed (such as the divisibility condition for the invariant factors or the condition that a be the power of a prime in the case of the elementary divisors). To decide whether two modules are isomorphic it is necessary to first write them in such a standard (or *canonical*) form.

EXERCISES

1. Let M be a module over the integral domain R .

- (a) Suppose x is a nonzero torsion element in M . Show that x and 0 are “linearly dependent.” Conclude that the rank of $\text{Tor}(M)$ is 0, so that in particular any torsion R -module has rank 0.
- (b) Show that the rank of M is the same as the rank of the (torsion free) quotient $M/\text{Tor}M$.

2. Let M be a module over the integral domain R .

- (a) Suppose that M has rank n and that x_1, x_2, \dots, x_n is any maximal set of linearly independent elements of M . Let $N = R x_1 + \dots + R x_n$ be the submodule generated by x_1, x_2, \dots, x_n . Prove that N is isomorphic to R^n and that the quotient M/N is a torsion R -module (equivalently, the elements x_1, \dots, x_n are linearly independent and for any $y \in M$ there is a nonzero element $r \in R$ such that ry can be written as a linear combination $r_1x_1 + \dots + r_nx_n$ of the x_i).
- (b) Prove conversely that if M contains a submodule N that is free of rank n (i.e., $N \cong R^n$) such that the quotient M/N is a torsion R -module then M has rank n . [Let y_1, y_2, \dots, y_{n+1} be any $n+1$ elements of M . Use the fact that M/N is torsion to write $r_i y_i$ as a linear combination of a basis for N for some nonzero elements r_1, \dots, r_{n+1} of R . Use an argument as in the proof of Proposition 3 to see that the $r_i y_i$, and hence also the y_i , are linearly dependent.]

3. Let R be an integral domain and let A and B be R -modules of ranks m and n , respectively. Prove that the rank of $A \oplus B$ is $m + n$. [Use the previous exercise.]
4. Let R be an integral domain, let M be an R -module and let N be a submodule of M . Suppose M has rank n , N has rank r and the quotient M/N has rank s . Prove that $n = r + s$. [Let x_1, x_2, \dots, x_s be elements of M whose images in M/N are a maximal set of independent elements and let $x_{s+1}, x_{s+2}, \dots, x_{s+r}$ be a maximal set of independent elements in N . Prove that x_1, x_2, \dots, x_{s+r} are linearly independent in M and that for any element $y \in M$ there is a nonzero element $r \in R$ such that ry is a linear combination of these elements. Then use Exercise 2.]
5. Let $R = \mathbb{Z}[x]$ and let $M = (2, x)$ be the ideal generated by 2 and x , considered as a submodule of R . Show that $(2, x)$ is not a basis of M . [Find a nontrivial R -linear dependence between these two elements.] Show that the rank of M is 1 but that M is not free of rank 1 (cf. Exercise 2).
6. Show that if R is an integral domain and M is any nonprincipal ideal of R then M is torsion free of rank 1 but is not a free R -module.

7. Let R be any ring, let A_1, A_2, \dots, A_m be R -modules and let B_i be a submodule of A_i , $1 \leq i \leq m$. Prove that

$$(A_1 \oplus A_2 \oplus \cdots \oplus A_m) / (B_1 \oplus B_2 \oplus \cdots \oplus B_m) \cong (A_1 / B_1) \oplus (A_2 / B_2) \oplus \cdots \oplus (A_m / B_m).$$

8. Let R be a P.I.D., let B be a torsion R -module and let p be a prime in R . Prove that if $pb = 0$ for some nonzero $b \in B$, then $\text{Ann}(B) \subseteq (p)$.

9. Give an example of an integral domain R and a nonzero torsion R -module M such that $\text{Ann}(M) = 0$. Prove that if N is a finitely generated torsion R -module then $\text{Ann}(N) \neq 0$.
10. For p a prime in the P.I.D. R and N an R -module prove that the p -primary component of N is a submodule of N and prove that N is the direct sum of its p -primary components (there need not be finitely many of them).

11. Let R be a P.I.D., let a be a nonzero element of R and let $M = R/(a)$. For any prime p of R prove that

$$p^{k-1}M/p^kM \cong \begin{cases} R/(p) & \text{if } k \leq n \\ 0 & \text{if } k > n, \end{cases}$$

where n is the power of p dividing a in R .

12. Let R be a P.I.D. and let p be a prime in R .

- (a) Let M be a finitely generated torsion R -module. Use the previous exercise to prove that $p^{k-1}M/p^kM \cong F^{n_k}$ where F is the field $R/(p)$ and n_k is the number of elementary divisors of M which are powers p^α with $\alpha \geq k$.
- (b) Suppose M_1 and M_2 are isomorphic finitely generated torsion R -modules. Use (a) to prove that, for every $k \geq 0$, M_1 and M_2 have the same number of elementary divisors p^α with $\alpha \geq k$. Prove that this implies M_1 and M_2 have the same set of elementary divisors.

13. If M is a finitely generated module over the P.I.D. R , describe the structure of $M/\text{Tor}(M)$.
14. Let R be a P.I.D. and let M be a torsion R -module. Prove that M is irreducible (cf. Exercises 9 to 11 of Section 10.3) if and only if $M = Rm$ for any nonzero element $m \in M$ where the annihilator of m is a nonzero prime ideal (p) .
15. Prove that if R is a Noetherian ring then R^n is a Noetherian R -module. [Fix a basis of R^n . If M is a submodule of R^n show that the collection of first coordinates of elements of M is a submodule of R hence is finitely generated. Let m_1, m_2, \dots, m_k be elements of M

whose first coordinates generate this submodule of R . Show that any element of M can be written as an R -linear combination of m_1, m_2, \dots, m_k plus an element of M whose first coordinate is 0. Prove that $M \cap R^{n-1}$ is a submodule of R^{n-1} where R^{n-1} is the set of elements of R^n with first coordinate 0 and then use induction on n .

The following set of exercises outlines a proof of Theorem 5 in the special case where R is a Euclidean Domain using a matrix argument involving row and column operations. This applies in particular to the cases $R = \mathbb{Z}$ and $R = F[x]$ of interest in the applications and is computationally useful.

Let R be a Euclidean Domain and let M be an R -module.

- 16.** Prove that M is finitely generated if and only if there is a surjective R -homomorphism $\varphi : R^n \rightarrow M$ for some integer n (this is true for any ring R).

Suppose $\varphi : R^n \rightarrow M$ is a surjective R -module homomorphism. By Exercise 15, $\ker \varphi$ is finitely generated. If x_1, x_2, \dots, x_n is a basis for R^n and y_1, \dots, y_m are generators for $\ker \varphi$ we have

$$y_i = a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n \quad i = 1, 2, \dots, m$$

with coefficients $a_{ij} \in R$. It follows that the homomorphism φ (hence the module structure of M) is determined by the choice of generators for R^n and the matrix $A = (a_{ij})$. Such a matrix A will be called a *relations matrix*.

- 17.** (a) Show that interchanging x_i and x_j in the basis for R^n interchanges the i^{th} column with the j^{th} column in the corresponding relations matrix.
 (b) Show that, for any $a \in R$, replacing the element x_j by $x_j - ax_i$ in the basis for R^n gives another basis for R^n and that the corresponding relations matrix for this basis is the same as the original relations matrix except that a times the j^{th} column has been added to the i^{th} column. [Note that $\cdots + a_i x_i + \cdots + a_j x_j + \cdots = \cdots + (a_i + aa_j)x_i + \cdots + a_j(x_j - ax_i) + \cdots$.]
- 18.** (a) Show that interchanging the generators y_i and y_j interchanges the i^{th} row with the j^{th} row in the relations matrix.
 (b) Show that, for any $a \in R$, replacing the element y_j by $y_j - ay_i$ gives another set of generators for $\ker \varphi$ and that the corresponding relations matrix for this choice of generators is the same as the original relations matrix except that $-a$ times the i^{th} row has been added to the j^{th} row.
- 19.** By the previous two exercises we may perform elementary row and column operations on a given relations matrix by choosing different generators for R^n and $\ker \varphi$. If all relation matrices are the zero matrix then $\ker \varphi = 0$ and $M \cong R^n$. Otherwise let a_1 be the (nonzero) g.c.d. (recall R is a Euclidean Domain) of all the entries in a fixed initial relations matrix for M .
 (a) Prove that by elementary row and column operations we may assume a_1 occurs in a relations matrix of the form

$$\begin{pmatrix} a_1 & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

where a_1 divides a_{ij} , $i = 1, 2, \dots, m$, $j = 1, 2, \dots, n$.

- (b) Prove that there is a relations matrix of the form

$$\begin{pmatrix} a_1 & 0 & \dots & 0 \\ 0 & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{m2} & \dots & a_{mn} \end{pmatrix}$$

where a_1 divides all the entries.

- (c) Let a_2 be a g.c.d. of all the entries except the element a_1 in the relations matrix in (b). Prove that there is a relations matrix of the form

$$\begin{pmatrix} a_1 & 0 & 0 & \dots & 0 \\ 0 & a_2 & 0 & \dots & 0 \\ 0 & 0 & a_{33} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & a_{m3} & \dots & a_{mn} \end{pmatrix}$$

where a_1 divides a_2 and a_2 divides all the other entries of the matrix.

- (d) Prove that there is a relations matrix of the form $\begin{pmatrix} D & 0 \\ 0 & 0 \end{pmatrix}$ where D is a diagonal matrix with nonzero entries a_1, a_2, \dots, a_k , $k \leq n$, satisfying

$$a_1 \mid a_2 \mid \dots \mid a_k.$$

Conclude that

$$M \cong R/(a_1) \oplus R/(a_2) \oplus \dots \oplus R/(a_k) \oplus R^{n-k}.$$

If n is not the minimal number of generators required for M then some of the initial elements a_1, a_2, \dots above will be units, so the corresponding direct summands above will be 0. If we remove these irrelevant factors we have produced the invariant factors of the module M . Further, the image of the new generators for R^n corresponding to the direct summands above will then be a set of R -generators for the cyclic submodules of M in its invariant factor decomposition (note that the image in M of the generators corresponding to factors with a_i a unit will be 0). The *column* operations performed in the relations matrix reduction correspond to changing the basis used for R^n as described in Exercise 17:

- (a) Interchanging the i^{th} column with the j^{th} column corresponds to interchanging the i^{th} and j^{th} elements in the basis for R^n .
- (b) For any $a \in R$, adding a times the j^{th} column to the i^{th} column corresponds to subtracting a times the i^{th} basis element from the j^{th} basis element.

Keeping track of the column operations performed and changing the initial choice of generators for M in the same way therefore gives a set of R -generators for the cyclic submodules of M in its invariant factor decomposition.

This process is quite fast computationally once an initial set of generators for M and initial relations matrix are determined. The element a_1 is determined using the Euclidean Algorithm as the g.c.d. of the elements in the initial relations matrix. Using the row and column operations we can obtain the appropriate linear combination of the entries to produce this g.c.d. in the (1,1)-position of a new relations matrix. One then subtracts the appropriate multiple of the first column and first row to obtain a matrix as in Exercise 19(b), then iterates this process. Some examples of this procedure in a special case are given at the end of the following section.

20. Let R be an integral domain with quotient field F and let M be any R -module. Prove that the rank of M equals the dimension of the vector space $F \otimes_R M$ over F .

21. Prove that a finitely generated module over a P.I.D. is projective if and only if it is free.
22. Let R be a P.I.D. that is not a field. Prove that no finitely generated R -module is injective. [Use Exercise 4, Section 10.5 to consider torsion and free modules separately.]

12.2 THE RATIONAL CANONICAL FORM

We now apply our results on finitely generated modules in the special case where the P.I.D. is the ring $F[x]$ of polynomials in x with coefficients in a field F .

Let V be a finite dimensional vector space over F of dimension n and let T be a fixed linear transformation of V (i.e., from V to itself). As we saw in Chapter 10 we can consider V as an $F[x]$ -module where the element x acts on V as the linear transformation T (and so any polynomial in x acts on V as the same polynomial in T). Since V has finite dimension over F by assumption, it is by definition finitely generated as an F -module, hence certainly finitely generated as an $F[x]$ -module, so the classification theorems of the preceding section apply.

Any nonzero free $F[x]$ -module (being isomorphic to a direct sum of copies of $F[x]$) is an infinite dimensional vector space over F , so if V has finite dimension over F then it must in fact be a torsion $F[x]$ -module (i.e., its free rank is 0). It follows from the Fundamental Theorem that then V is isomorphic as an $F[x]$ -module to the direct sum of cyclic, torsion $F[x]$ -modules. We shall see that this decomposition of V will allow us to choose a basis for V with respect to which the matrix representation for the linear transformation T is in a specific simple form. When we use the invariant factor decomposition of V we obtain the *rational canonical form* for the matrix for T , which we analyze in this section. When we use the elementary divisor decomposition (and when F contains all the eigenvalues of T) we obtain the *Jordan canonical form*, considered in the following section and mentioned earlier as the matrix representing T which is as close to being a diagonal matrix as possible. The uniqueness portion of the Fundamental Theorem ensures that the rational and Jordan canonical forms are unique (which is why they are referred to as *canonical*).

One important use of these canonical forms is to classify the distinct linear transformations of V . In particular they allow us to determine when two matrices represent the same linear transformation, i.e., when two given $n \times n$ matrices are similar.

Note that this will be another instance where the structure of the space being acted upon (the invariant factor decomposition of V for example) is used to obtain significant information on the algebraic objects (in this case the linear transformations) which are acting. This will be considered in the case of *groups* acting on vector spaces in Chapter 18 (and goes under the name of Representation Theory of Groups).

Before describing the rational canonical form in detail we first introduce some linear algebra.

Definition.

- (1) An element λ of F is called an *eigenvalue* of the linear transformation T if there is a nonzero vector $v \in V$ such that $T(v) = \lambda v$. In this situation v is called an *eigenvector* of T with corresponding eigenvalue λ .

- (2) If A is an $n \times n$ matrix with coefficients in F , an element λ is called an *eigenvalue* of A with corresponding eigenvector v if v is a nonzero $n \times 1$ column vector such that $Av = \lambda v$.
- (3) If λ is an eigenvalue of the linear transformation T , the set $\{v \in V \mid T(v) = \lambda v\}$ is called the *eigenspace* of T corresponding to the eigenvalue λ . Similarly, if λ is an eigenvalue of the $n \times n$ matrix A , the set of $n \times 1$ matrices v with $Av = \lambda v$ is called the *eigenspace* of A corresponding to the eigenvalue λ .

Note that if we fix a basis \mathcal{B} of V then any linear transformation T of V has an associated $n \times n$ matrix A . Conversely, if A is any $n \times n$ matrix then the map T defined by $T(v) = Av$ for $v \in V$, where the v on the right is the $n \times 1$ vector consisting of the coordinates of v with respect to the fixed basis \mathcal{B} of V , is a linear transformation of V . Then v is an eigenvector of T with corresponding eigenvalue λ if and only if the coordinate vector of v with respect to \mathcal{B} is an eigenvector of A with eigenvalue λ . In other words, the eigenvalues for the linear transformation T are the same as the eigenvalues for the matrix A of T with respect to any fixed basis for V .

Definition. The determinant of a linear transformation from V to V is the determinant of any matrix representing the linear transformation (note that this does not depend on the choice of the basis used).

Proposition 12. The following are equivalent:

- (1) λ is an eigenvalue of T
- (2) $\lambda I - T$ is a singular linear transformation of V
- (3) $\det(\lambda I - T) = 0$.

Proof: Since λ is an eigenvalue of T with corresponding eigenvector v if and only if v is a nonzero vector in the kernel of $\lambda I - T$, it follows that (1) and (2) are equivalent.

(2) and (3) are equivalent by our results on determinants.

Definition. Let x be an indeterminate over F . The polynomial $\det(xI - T)$ is called the *characteristic polynomial* of T and will be denoted $c_T(x)$. If A is an $n \times n$ matrix with coefficients in F , $\det(xI - A)$ is called the *characteristic polynomial* of A and will be denoted $c_A(x)$.

It is easy to see by expanding the determinant that the characteristic polynomial of either T or A is a monic polynomial of degree $n = \dim V$. Proposition 12 says that the set of eigenvalues of T (or A) is precisely the set of roots of the characteristic polynomial of T (of A , respectively). In particular, T has at most n distinct eigenvalues.

We have seen that V considered as a module over $F[x]$ via the linear transformation T is a torsion $F[x]$ -module. Let $m(x) \in F[x]$ be the unique monic polynomial generating the annihilator of V in $F[x]$. Equivalently, $m(x)$ is the unique monic polynomial of minimal degree annihilating V (i.e., such that $m(T)$ is the 0 linear transformation), and if $f(x) \in F[x]$ is any polynomial annihilating V , $m(x)$ divides $f(x)$. Since the ring of all $n \times n$ matrices over F is isomorphic to the collection of all linear transformations of V to itself (an isomorphism is obtained by choosing a basis for V), it follows that for

any $n \times n$ matrix A over F there is similarly a unique monic polynomial of minimal degree with $m(A)$ the zero matrix.

Definition. The unique monic polynomial which generates the ideal $\text{Ann}(V)$ in $F[x]$ is called the *minimal polynomial* of T and will be denoted $m_T(x)$. The unique monic polynomial of smallest degree which when evaluated at the matrix A is the zero matrix is called the *minimal polynomial* of A and will be denoted $m_A(x)$.

It is easy to see (cf. Exercise 5) that the degrees of these minimal polynomials are at most n^2 where n is the dimension of V . We shall shortly prove that the minimal polynomial for T is a divisor of the characteristic polynomial for T (this is the *Cayley-Hamilton Theorem*), and similarly for A , so in fact the degrees of these polynomials are at most n .

We now describe the *rational canonical form* of the linear transformation T (respectively, of the $n \times n$ matrix A). By Theorem 5 we have an isomorphism

$$V \cong F[x]/(a_1(x)) \oplus F[x]/(a_2(x)) \oplus \cdots \oplus F[x]/(a_m(x)) \quad (12.1)$$

of $F[x]$ -modules where $a_1(x), a_2(x), \dots, a_m(x)$ are polynomials in $F[x]$ of degree at least one with the divisibility conditions

$$a_1(x) \mid a_2(x) \mid \cdots \mid a_m(x).$$

These invariant factors $a_i(x)$ are only determined up to a unit in $F[x]$ but since the units of $F[x]$ are precisely the nonzero elements of F (i.e., the nonzero constant polynomials), we may make these polynomials *unique* by stipulating that they be *monic*.

Since the annihilator of V is the ideal $(a_m(x))$ (part (3) of Theorem 5), we immediately obtain:

Proposition 13. The minimal polynomial $m_T(x)$ is the largest invariant factor of V . All the invariant factors of V divide $m_T(x)$.

We shall see below how to calculate not only the minimal polynomial for T but also the other invariant factors.

We now choose a basis for each of the direct summands for V in the decomposition (1) above for which the matrix for T is quite simple. Recall that the linear transformation T acting on the left side of (1) is the element x acting by multiplication on each of the factors on the right side of the isomorphism in (1).

We have seen in the example following Proposition 1 of Chapter 11 that the elements $1, \bar{x}, \bar{x}^2, \dots, \bar{x}^{k-1}$ give a basis for the vector space $F[x]/(a(x))$ where $a(x) = x^k + b_{k-1}x^{k-1} + \cdots + b_1x + b_0$ is any monic polynomial in $F[x]$ and $\bar{x} = x \bmod (a(x))$. With respect to this basis the linear transformation of multiplication by x acts in a simple manner:

$$\begin{array}{rcl} & 1 & \mapsto \bar{x} \\ & \bar{x} & \mapsto \bar{x}^2 \\ & \bar{x}^2 & \mapsto \bar{x}^3 \\ x : & \vdots & \\ & \bar{x}^{k-2} & \mapsto \bar{x}^{k-1} \\ & \bar{x}^{k-1} & \mapsto \bar{x}^k = -b_0 - b_1\bar{x} - \cdots - b_{k-1}\bar{x}^{k-1} \end{array}$$

where the last equality is because $\bar{x}^k + b_{k-1}\bar{x}^{k-1} + \cdots + b_1\bar{x} + b_0 = 0$ since $a(\bar{x}) = 0$ in $F[x]/(a(x))$. With respect to this basis, the matrix for multiplication by x is therefore

$$\begin{pmatrix} 0 & 0 & \dots & \dots & \dots & -b_0 \\ 1 & 0 & \dots & \dots & \dots & -b_1 \\ 0 & 1 & \dots & \dots & \dots & -b_2 \\ 0 & 0 & \ddots & & & \vdots \\ \vdots & \vdots & & \ddots & & \vdots \\ 0 & 0 & \dots & \dots & 1 & -b_{k-1} \end{pmatrix}.$$

Such matrices are given a name:

Definition. Let $a(x) = x^k + b_{k-1}x^{k-1} + \cdots + b_1x + b_0$ be any monic polynomial in $F[x]$. The *companion matrix* of $a(x)$ is the $k \times k$ matrix with 1's down the first subdiagonal, $-b_0, -b_1, \dots, -b_{k-1}$ down the last column and zeros elsewhere. The companion matrix of $a(x)$ will be denoted by $C_{a(x)}$.

We apply this to each of the cyclic modules on the right side of (1) above and let B_i be the elements of V corresponding to the basis chosen above for the cyclic factor $F[x]/(a_i(x))$ under the isomorphism in (1). Then by definition the linear transformation T acts on B_i by the companion matrix for $a_i(x)$ since we have seen that this is how multiplication by x acts. The union B of the B_i 's gives a basis for V since the sum on the right of (1) is direct and with respect to this basis the linear transformation T has as matrix the *direct sum* of the companion matrices for the invariant factors, i.e.,

$$\begin{pmatrix} C_{a_1(x)} & & & \\ & C_{a_2(x)} & & \\ & & \ddots & \\ & & & C_{a_m(x)} \end{pmatrix}. \quad (12.2)$$

Notice that this matrix is uniquely determined from the invariant factors of the $F[x]$ -module V and, by Theorem 9, the list of invariant factors uniquely determines the module V up to isomorphism as an $F[x]$ -module.

Definition.

- (1) A matrix is said to be in *rational canonical form* if it is the direct sum of companion matrices for monic polynomials $a_1(x), \dots, a_m(x)$ of degree at least one with $a_1(x) \mid a_2(x) \mid \cdots \mid a_m(x)$. The polynomials $a_i(x)$ are called the *invariant factors* of the matrix. Such a matrix is also said to be a *block diagonal* matrix with blocks the companion matrices for the $a_i(x)$.
- (2) A *rational canonical form* for a linear transformation T is a matrix representing T which is in rational canonical form.

We have seen that any linear transformation T has a rational canonical form. We now see that this rational canonical form is unique (hence is called *the* rational canonical form for T). To see this note that the process we used to determine the matrix of T

from the direct sum decomposition is reversible. Suppose $b_1(x), b_2(x), \dots, b_t(x)$ are monic polynomials in $F[x]$ of degree at least one such that $b_i(x) \mid b_{i+1}(x)$ for all i and suppose for some basis \mathcal{E} of V , that the matrix of T with respect to the basis \mathcal{E} is the direct sum of the companion matrices of the $b_i(x)$. Then V must be a direct sum of T -stable subspaces D_i , one for each $b_i(x)$ in such a way that the matrix of T on each D_i is the companion matrix of $b_i(x)$. Let \mathcal{E}_i be the corresponding (ordered) basis of D_i (so \mathcal{E} is the union of the \mathcal{E}_i) and let e_i be the first basis element in \mathcal{E}_i . Then it is easy to see that D_i is a cyclic $F[x]$ -module with generator e_i and that the annihilator of D_i is $b_i(x)$. Thus the torsion $F[x]$ -module V decomposes into a direct sum of cyclic $F[x]$ -modules in two ways, both of which satisfy the conditions of Theorem 5, i.e., both of which give lists of invariant factors. Since the invariant factors are unique by Theorem 9, $a_i(x)$ and $b_i(x)$ must differ by a unit factor in $F[x]$ and since the polynomials are monic by assumption, we must have $a_i(x) = b_i(x)$ for all i . This proves the following result:

Theorem 14. (*Rational Canonical Form for Linear Transformations*) Let V be a finite dimensional vector space over the field F and let T be a linear transformation of V .

- (1) There is a basis for V with respect to which the matrix for T is in rational canonical form, i.e., is a block diagonal matrix whose diagonal blocks are the companion matrices for monic polynomials $a_1(x), a_2(x), \dots, a_m(x)$ of degree at least one with $a_1(x) \mid a_2(x) \mid \dots \mid a_m(x)$.
- (2) The rational canonical form for T is unique.

The use of the word *rational* is to indicate that this canonical form is calculated entirely within the field F and exists for any linear transformation T . This is not the case for the Jordan canonical form (considered later), which only exists if the field F contains the eigenvalues for T (cf. also the remarks following Corollary 18).

The following result translates the notion of similar linear transformations (i.e., the same linear transformation up to a change of basis) into the language of modules and relates this notion to rational canonical forms.

Theorem 15. Let S and T be linear transformations of V . Then the following are equivalent:

- (1) S and T are similar linear transformations
- (2) the $F[x]$ -modules obtained from V via S and via T are isomorphic $F[x]$ -modules
- (3) S and T have the same rational canonical form.

Proof: [(1) implies (2)] Assume there is a nonsingular linear transformation U such that $S = UTU^{-1}$. The vector space isomorphism $U : V \rightarrow V$ is also an $F[x]$ -module homomorphism, where x acts on the first V via T and on the second via S , since for example $U(xv) = U(Tv) = UT(v) = SU(v) = x(Uv)$. Hence this is an $F[x]$ -module isomorphism of the two modules in (2).

[(2) implies (3)] Assume (2) holds and denote by V_1 the vector space V made into an $F[x]$ -module via S and denote by V_2 the space V made into an $F[x]$ -module via T . Since $V_1 \cong V_2$ as $F[x]$ -modules they have the same list of invariant factors. Thus S and T have a common rational canonical form.

[(3) implies (1)] Assume (3) holds. Since S and T have the same matrix representation with respect to some choice of (possibly different) bases of V by assumption, they are, up to a change of basis, the same linear transformation of V , hence are similar.

Let A be any $n \times n$ matrix with entries from F . Let V be an n -dimensional vector space over F . Recall we can then *define* a linear transformation T on V by choosing a basis for V and setting $T(v) = Av$ where v on the right hand side means the $n \times 1$ column vector of coordinates of v with respect to our chosen basis (this is just the usual identification of linear transformations with matrices). Then (of course) the matrix for this T with respect to this basis is the given matrix A . Put another way, any $n \times n$ matrix A with entries from the field F arises as the matrix for some linear transformation T of an n -dimensional vector space.

This dictionary between linear transformations of vector spaces and matrices allows us to state our previous two results in the language of matrices:

Theorem 16. (*Rational Canonical Form for Matrices*) Let A be an $n \times n$ matrix over the field F .

- (1) The matrix A is similar to a matrix in rational canonical form, i.e., there is an invertible $n \times n$ matrix P over F such that $P^{-1}AP$ is a block diagonal matrix whose diagonal blocks are the companion matrices for monic polynomials $a_1(x), a_2(x), \dots, a_m(x)$ of degree at least one with $a_1(x) | a_2(x) | \dots | a_m(x)$.
- (2) The rational canonical form for A is unique.

Definition. The *invariant factors* of an $n \times n$ matrix over a field F are the invariant factors of its rational canonical form.

Theorem 17. Let A and B be $n \times n$ matrices over the field F . Then A and B are similar if and only if A and B have the same rational canonical form.

If A is a matrix with entries from a field F and F is a subfield of a larger field K then we may also consider A as a matrix over K . The next result shows that the rational canonical form for A and questions of similarity do not depend on which field contains the entries of A .

Corollary 18. Let A and B be two $n \times n$ matrices over a field F and suppose F is a subfield of the field K .

- (1) The rational canonical form of A is the same whether it is computed over K or over F . The minimal and characteristic polynomials and the invariant factors of A are the same whether A is considered as a matrix over F or as a matrix over K .
- (2) The matrices A and B are similar over K if and only if they are similar over F , i.e., there exists an invertible $n \times n$ matrix P with entries from K such that $B = P^{-1}AP$ if and only if there exists an (in general different) invertible $n \times n$ matrix Q with entries from F such that $B = Q^{-1}AQ$.

Proof: (1) Let M be the rational canonical form of A when computed over the smaller field F . Since M satisfies the conditions in the definition of the rational canonical form over K , the uniqueness of the rational canonical form implies that M is also

the rational canonical form of A over K . Hence the invariant factors of A are the same whether A is viewed over F or over K . In particular, since the minimal polynomial is the largest invariant factor of A it also does not depend on the field over which A is viewed. It is clear from the determinant definition of the characteristic polynomial of A that this polynomial depends only on the entries of A (we shall see shortly that the characteristic polynomial is the product of all the invariant factors for A , which will give an alternate proof of this result).

(2) If A and B are similar over the smaller field F they are clearly similar over K . Conversely, if A and B are similar over K , they have the same rational canonical form over K . By (1) they have the same rational canonical form over F , hence are similar over F by Theorem 17.

This corollary asserts in particular that the rational canonical form for an $n \times n$ matrix A is an $n \times n$ matrix with entries in the smallest field containing the entries of A . Further, this canonical form is the same matrix even if we allow conjugation of A by nonsingular matrices whose entries come from larger fields. This explains the terminology of *rational* canonical form.

The next proposition gives the connection between the characteristic polynomial of a matrix (or of a linear transformation) and its invariant factors and is quite useful for determining these invariant factors (particularly for matrices of small size).

Lemma 19. Let $a(x) \in F[x]$ be any monic polynomial.

- (1) The characteristic polynomial of the companion matrix of $a(x)$ is $a(x)$.
- (2) If M is the block diagonal matrix

$$M = \begin{pmatrix} A_1 & 0 & \dots & 0 \\ 0 & A_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & A_k \end{pmatrix},$$

given by the direct sum of matrices A_1, A_2, \dots, A_k then the characteristic polynomial of M is the product of the characteristic polynomials of A_1, A_2, \dots, A_k .

Proof: These are both straightforward exercises.

Proposition 20. Let A be an $n \times n$ matrix over the field F .

- (1) The characteristic polynomial of A is the product of all the invariant factors of A .
- (2) (*The Cayley–Hamilton Theorem*) The minimal polynomial of A divides the characteristic polynomial of A .
- (3) The characteristic polynomial of A divides some power of the minimal polynomial of A . In particular these polynomials have the same roots, not counting multiplicities.

The same statements are true if the matrix A is replaced by a linear transformation T of an n -dimensional vector space over F .

Proof: Let B be the rational canonical form of A . By the previous lemma the block diagonal form of B shows that the characteristic polynomial of B is the product of the characteristic polynomials of the companion matrices of the invariant factors of A . By the first part of the lemma above, the characteristic polynomial of the companion matrix $C_{a(x)}$ for $a(x)$ is just $a(x)$, which implies that the characteristic polynomial for B is the product of the invariant factors of A . Since A and B are similar, they have the same characteristic polynomial, which proves (1). Assertion (2) is immediate from (1) since the minimal polynomial for A is the largest invariant factor of A . The fact that all the invariant factors divide the largest one immediately implies (3). The final assertion is clear from the dictionary between linear transformations of vector spaces and matrices.

Note that part (2) of the proposition is the assertion that the matrix A satisfies its own characteristic polynomial, i.e., $c_A(A) = 0$ as matrices, which is the usual formulation for the Cayley–Hamilton Theorem. Note also that it implies the degree of the minimal polynomial for A has degree at most n , a result mentioned before.

The relations in Proposition 20 are frequently quite useful in the determination of the invariant factors for a matrix A , particularly for matrices of small degree (cf. Exercises 3 and 4 and the examples). The following result (which relies on Exercises 16 to 19 in the previous section and whose proof we outline in the exercises) computes the invariant factors in general.

Let A be an $n \times n$ matrix over the field F . Then $xI - A$ is an $n \times n$ matrix with entries in $F[x]$. The three operations

- (a) interchanging two rows or columns
 - (b) adding a multiple (in $F[x]$) of one row or column to another
 - (c) multiplying any row or column by a unit in $F[x]$, i.e., by a nonzero element in F ,
- are called *elementary row and column operations*.

Theorem 21. Let A be an $n \times n$ matrix over the field F . Using the three elementary row and column operations above, the $n \times n$ matrix $xI - A$ with entries from $F[x]$ can be put into the diagonal form (called the *Smith Normal Form* for A)

$$\begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & a_1(x) & \\ & & & & a_2(x) \\ & & & & \ddots \\ & & & & a_m(x) \end{pmatrix}$$

with monic nonzero elements $a_1(x), a_2(x), \dots, a_m(x)$ of $F[x]$ with degrees at least one and satisfying $a_1(x) | a_2(x) | \dots | a_m(x)$. The elements $a_1(x), \dots, a_m(x)$ are the invariant factors of A .

Proof: cf. the exercises.

Invariant Factor Decomposition Algorithm: Converting to Rational Canonical Form

As mentioned in the exercises near the end of the previous section, keeping track of the operations necessary to diagonalize $xI - A$ will explicitly give a matrix P such that $P^{-1}AP$ is in rational canonical form. Equivalently, if V is a given $F[x]$ -module with vector space basis $[e_1, e_2, \dots, e_n]$, then P defines the change of basis giving the Invariant Factor Decomposition of V into a direct sum of cyclic $F[x]$ -modules. In particular, if A is the matrix of the linear transformation T of the $F[x]$ -module V defined by x (i.e., $T(e_j) = xe_j = \sum_{i=1}^n a_{ij}e_i$ where $A = (a_{ij})$), then the matrix P defines the change of basis for V with respect to which the matrix for T is in rational canonical form.

We first describe the algorithm in the general context of determining the Invariant Factor Decomposition of a given $F[x]$ -module V with vector space basis $[e_1, e_2, \dots, e_n]$ (the proof is outlined in the exercises). We then describe the algorithm to convert a given $n \times n$ matrix A to rational canonical form (in which reference to an underlying vector space and associated linear transformation are suppressed).

Explicit numerical examples of this algorithm are given in Examples 2 and 3 following.

Invariant Factor Decomposition Algorithm

Let V be an $F[x]$ -module with vector space basis $[e_1, e_2, \dots, e_n]$ (so in particular these elements are generators for V as an $F[x]$ -module). Let T be the linear transformation of V to itself defined by x and let A be the $n \times n$ matrix associated to T and this choice of basis for V , i.e.,

$$T(e_j) = xe_j = \sum_{i=1}^n a_{ij}e_i \quad \text{where} \quad A = (a_{ij}).$$

- (1) Use the following three elementary row and column operations to diagonalize the matrix $xI - A$ over $F[x]$, keeping track of the *row* operations used:
 - (a) interchange two rows or columns (which will be denoted by $R_i \leftrightarrow R_j$ for the interchange of the i^{th} and j^{th} rows and similarly by $C_i \leftrightarrow C_j$ for columns),
 - (b) add a multiple (in $F[x]$) of one row or column to another (which will be denoted by $R_i + p(x)R_j \mapsto R_i$ if $p(x)$ times the j^{th} row is added to the i^{th} row, and similarly by $C_i + p(x)C_j \mapsto C_i$ for columns),
 - (c) multiply any row or column by a unit in $F[x]$, i.e., by a nonzero element in F (which will be denoted by uR_i if the i^{th} row is multiplied by $u \in F^\times$, and similarly by uC_i for columns).
- (2) Beginning with the $F[x]$ -module generators $[e_1, e_2, \dots, e_n]$, for each row operation used in (1), change the set of generators by the following rules:
 - (a) If the i^{th} row is interchanged with the j^{th} row then interchange the i^{th} and j^{th} generators.
 - (b) If $p(x)$ times the j^{th} row is added to the i^{th} row then subtract $p(x)$ times the i^{th} generator from the j^{th} generator (note the indices).

- (c) If the i^{th} row is multiplied by the unit $u \in F$ then divide the i^{th} generator by u .
- (3) When $xI - A$ has been diagonalized to the form in Theorem 21 the generators $[e_1, e_2, \dots, e_n]$ for V will be in the form of $F[x]$ -linear combinations of e_1, e_2, \dots, e_n . Use $xe_j = T(e_j) = \sum_{i=1}^n a_{ij}e_i$ to write these elements as F -linear combinations of e_1, e_2, \dots, e_n . When $xI - A$ has been diagonalized, the first $n - m$ of these linear combinations are 0 (providing a useful numerical check on the computations) and the remaining m linear combinations are nonzero, i.e., the generators for V are in the form $[0, \dots, 0, f_1, \dots, f_m]$ corresponding precisely to the diagonal elements in Theorem 21. The elements f_1, \dots, f_m are a set of $F[x]$ -module generators for the cyclic factors in the invariant factor decomposition of V (with annihilators $(a_1(x)), \dots, (a_m(x))$, respectively):

$$V = F[x] f_1 \oplus F[x] f_2 \oplus \dots \oplus F[x] f_m,$$

$$F[x] f_i \cong F[x]/(a_i(x)) \quad i = 1, 2, \dots, m,$$

giving the Invariant Factor Decomposition of the $F[x]$ -module V .

- (4) The corresponding *vector space* basis for each cyclic factor of V is then given by the elements $f_i, Tf_i, T^2 f_i, \dots, T^{\deg a_i(x)-1} f_i$.
- (5) Write the k^{th} element of the vector space basis computed in (4) in terms of the original vector space basis $[e_1, e_2, \dots, e_n]$ and use the coordinates for the k^{th} column of an $n \times n$ matrix P . Then $P^{-1}AP$ is in rational canonical form (with diagonal blocks the companion matrices for the $a_i(x)$). This is the matrix for the linear transformation T with respect to the vector space basis in (4).

We now describe the algorithm to convert a given $n \times n$ matrix A to rational canonical form, i.e., to determine an $n \times n$ matrix P so that $P^{-1}AP$ is in rational canonical form. This is nothing more than the algorithm above applied to the vector space $V = F^n$ of $n \times 1$ column vectors with standard basis $[e_1, e_2, \dots, e_n]$ (where e_i is the column vector with 1 in the i^{th} position and 0's elsewhere) and T is the linear transformation defined by A and this choice of basis. Explicit reference to this underlying vector space and associated linear transformation are suppressed, so the algorithm is purely matrix theoretic.

Converting an $n \times n$ Matrix to Rational Canonical Form

Let A be an $n \times n$ matrix with entries in the field F .

- (1) Use the following three elementary row and column operations to diagonalize the matrix $xI - A$ over $F[x]$, keeping track of the *row* operations used:
- (a) interchange two rows or columns (which will be denoted by $R_i \leftrightarrow R_j$ for the interchange of the i^{th} and j^{th} rows and similarly by $C_i \leftrightarrow C_j$ for columns),
 - (b) add a multiple (in $F[x]$) of one row or column to another (which will be denoted by $R_i + p(x)R_j \mapsto R_i$ if $p(x)$ times the j^{th} row is added to the i^{th} row, and similarly by $C_i + p(x)C_j \mapsto C_i$ for columns),
 - (c) multiply any row or column by a unit in $F[x]$, i.e., by a nonzero element in F (which will be denoted by uR_i if the i^{th} row is multiplied by $u \in F^\times$, and similarly by uC_i for columns).

Define d_1, \dots, d_m to be the degrees of the monic nonconstant polynomials $a_1(x), \dots, a_m(x)$ appearing on the diagonal, respectively.

- (2) Beginning with the $n \times n$ identity matrix P' , for each row operation used in (1), change the matrix P' by the following rules:
 - (a) If $R_i \leftrightarrow R_j$ then interchange the i^{th} and j^{th} columns of P' (i.e., $C_i \leftrightarrow C_j$ for P').
 - (b) If $R_i + p(x)R_j \mapsto R_i$ then subtract the product of the matrix $p(A)$ times the i^{th} column of P' from the j^{th} column of P' (i.e., $C_j - p(A)C_i \mapsto C_j$ for P' — note the indices).
 - (c) If uR_i then divide the elements of the i^{th} column of P' by u (i.e., $u^{-1}C_i \mapsto C_i$ for P').
- (3) When $xI - A$ has been diagonalized to the form in Theorem 21 the first $n - m$ columns of the matrix P' are 0 (providing a useful numerical check on the computations) and the remaining m columns of P' are nonzero. For each $i = 1, 2, \dots, m$, multiply the i^{th} nonzero column of P' successively by $A^0 = I, A^1, A^2, \dots, A^{d_i-1}$, where d_i is the integer in (1) above and use the resulting column vectors (in this order) as the next d_i columns of an $n \times n$ matrix P . Then $P^{-1}AP$ is in rational canonical form (whose diagonal blocks are the companion matrices for the polynomials $a_1(x), \dots, a_m(x)$ in (1)).

In the theory of canonical forms for linear transformations (or matrices) the characteristic polynomial plays the role of the order of a finite abelian group and the minimal polynomial plays the role of the exponent (after all, they are the same invariants, one for modules over the Principal Ideal Domain \mathbb{Z} and the other for modules over the Principal Ideal Domain $F[x]$) so we can solve problems directly analogous to those we considered for finite abelian groups in Chapter 5. In particular, this includes the following:

- (A) determine the rational canonical form of a given matrix (analogous to decomposing a finite abelian group as a direct product of cyclic groups)
- (B) determine whether two given matrices are similar (analogous to determining whether two given finite abelian groups are isomorphic)
- (C) determine all similarity classes of matrices over F with a given characteristic polynomial (analogous to determining all abelian groups of a given order)
- (D) determine all similarity classes of $n \times n$ matrices over F with a given minimal polynomial (analogous to determining all abelian groups of rank at most n of a given exponent).

Examples

- (1) We find the rational canonical forms of the following matrices over \mathbb{Q} and determine if they are similar:

$$A = \begin{pmatrix} 2 & -2 & 14 \\ 0 & 3 & -7 \\ 0 & 0 & 2 \end{pmatrix} \quad B = \begin{pmatrix} 0 & -4 & 85 \\ 1 & 4 & -30 \\ 0 & 0 & 3 \end{pmatrix} \quad C = \begin{pmatrix} 2 & 2 & 1 \\ 0 & 2 & -1 \\ 0 & 0 & 3 \end{pmatrix}.$$

A direct computation shows that all three of these matrices have the same characteristic polynomial: $c_A(x) = c_B(x) = c_C(x) = (x - 2)^2(x - 3)$. Since the minimal and char-

acteristic polynomials have the same roots, the only possibilities for the minimal polynomials are $(x-2)(x-3)$ or $(x-2)^2(x-3)$. We quickly find that $(A-2I)(A-3I) = 0$, $(B-2I)(B-3I) \neq 0$ (the 1,1-entry is nonzero) and $(C-2I)(C-3I) \neq 0$ (the 1,2-entry is nonzero). It follows that

$$m_A(x) = (x-2)(x-3), \quad m_B(x) = m_C(x) = (x-2)^2(x-3).$$

It follows immediately that there are no additional invariant factors for B and C . Since the invariant factors for A divide the minimal polynomial and have product the characteristic polynomial, we see that A has for invariant factors the polynomials $x-2$, $(x-2)(x-3) = x^2 - 5x + 6$. (For 2×2 and 3×3 matrices the determination of the characteristic and minimal polynomials determines all the invariant factors, cf. Exercises 3 and 4.) We conclude that B and C are similar and neither is similar to A . The rational canonical forms are (note $(x-2)^2(x-3) = x^3 - 7x^2 + 16x - 12$)

$$\begin{pmatrix} 2 & 0 & 0 \\ 0 & 0 & -6 \\ 0 & 1 & 5 \end{pmatrix} \quad \begin{pmatrix} 0 & 0 & 12 \\ 1 & 0 & -16 \\ 0 & 1 & 7 \end{pmatrix} \quad \begin{pmatrix} 0 & 0 & 12 \\ 1 & 0 & -16 \\ 0 & 1 & 7 \end{pmatrix}.$$

- (2) In the example above the rational canonical forms were obtained simply by determining the characteristic and minimal polynomials for the matrices. As mentioned, this is sufficient for 2×2 and 3×3 matrices since this information is sufficient to determine all of the invariant factors. For larger matrices, however, this is in general not sufficient (cf. the next example) and more work is required to determine the invariant factors. In this example we again compute the rational canonical form for the matrix A in Example 1 following the two algorithms outlined above. While this is computationally more difficult for this small matrix (as will be apparent), it has the advantage even in this case that it also explicitly computes a matrix P with $P^{-1}AP$ in rational canonical form.

I. (*Invariant Factor Decomposition*) We use row and column operations (in $\mathbb{Q}[x]$) to reduce the matrix

$$xI - A = \begin{pmatrix} x-2 & 2 & -14 \\ 0 & x-3 & 7 \\ 0 & 0 & x-2 \end{pmatrix}$$

to diagonal form. As in the invariant factor decomposition algorithm, we shall use the notation $R_i \leftrightarrow R_j$ to denote the interchange of the i^{th} and j^{th} rows, $R_i + aR_j \mapsto R_i$ if a times the j^{th} row is added to the i^{th} row, simply uR_i if the i^{th} row is multiplied by u (and similarly for columns, using C instead of R). Note also that the first two operations we perform below are rather *ad hoc* and were chosen simply to have integers everywhere in the computation:

$$\begin{pmatrix} x-2 & 2 & -14 \\ 0 & x-3 & 7 \\ 0 & 0 & x-2 \end{pmatrix} \xrightarrow[R_1+R_2 \mapsto R_1]{\quad} \begin{pmatrix} x-2 & x-1 & -7 \\ 0 & x-3 & 7 \\ 0 & 0 & x-2 \end{pmatrix} \rightarrow$$

$$\xrightarrow[C_1-C_2 \mapsto C_1]{\quad} \begin{pmatrix} -1 & x-1 & -7 \\ -x+3 & x-3 & 7 \\ 0 & 0 & x-2 \end{pmatrix} \xrightarrow{-R_1} \begin{pmatrix} 1 & -x+1 & 7 \\ -x+3 & x-3 & 7 \\ 0 & 0 & x-2 \end{pmatrix} \rightarrow$$

$$\begin{array}{l}
R_2 + (x-3)R_1 \xrightarrow[\leftrightarrow R_2]{} \left(\begin{array}{ccc} 1 & -x+1 & 7 \\ 0 & -x^2+5x-6 & 7(x-2) \\ 0 & 0 & x-2 \end{array} \right) \xrightarrow[C_2+(x-1)C_1]{} \left(\begin{array}{ccc} 1 & 0 & 7 \\ 0 & -x^2+5x-6 & 7(x-2) \\ 0 & 0 & x-2 \end{array} \right) \rightarrow \\
\rightarrow_{C_3-7C_1} \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & -x^2+5x-6 & 7(x-2) \\ 0 & 0 & x-2 \end{array} \right) \xrightarrow[-C_2]{} \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & x^2-5x+6 & 7(x-2) \\ 0 & 0 & x-2 \end{array} \right) \rightarrow \\
\rightarrow_{R_2-7R_3} \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & x^2-5x+6 & 0 \\ 0 & 0 & x-2 \end{array} \right) \xrightarrow[R_2 \leftrightarrow R_3]{C_2 \leftrightarrow C_3} \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & x-2 & 0 \\ 0 & 0 & x^2-5x+6 \end{array} \right).
\end{array}$$

This determines the invariant factors $x - 2, x^2 - 5x + 6$ for this matrix, which we determined in Example 1 above. Let now V be a 3-dimensional vector space over \mathbb{Q} with basis e_1, e_2, e_3 and let T be the corresponding linear transformation (which defines the action of x on V), i.e.,

$$\begin{aligned}
xe_1 &= T(e_1) = 2e_1 \\
xe_2 &= T(e_2) = -2e_1 + 3e_2 \\
xe_3 &= T(e_3) = 14e_1 - 7e_2 + 2e_3.
\end{aligned}$$

The row operations used in the reduction above were

$$R_1 + R_2 \mapsto R_1, \quad -R_1, \quad R_2 + (x-3)R_1 \mapsto R_2, \quad R_2 - 7R_3 \mapsto R_2, \quad R_2 \leftrightarrow R_3.$$

Starting with the basis $[e_1, e_2, e_3]$ for V and changing it according to the rules given in the text, we obtain

$$\begin{aligned}
[e_1, e_2, e_3] &\longrightarrow [e_1, e_2 - e_1, e_3] \longrightarrow [-e_1, e_2 - e_1, e_3] \\
&\longrightarrow [-e_1 - (x-3)(e_2 - e_1), e_2 - e_1, e_3] \\
&\longrightarrow [-e_1 - (x-3)(e_2 - e_1), e_2 - e_1, e_3 + 7(e_2 - e_1)] \\
&\longrightarrow [-e_1 - (x-3)(e_2 - e_1), e_3 + 7(e_2 - e_1), e_2 - e_1].
\end{aligned}$$

Using the formulas above for the action of x , we see that these last elements are the elements $[0, -7e_1 + 7e_2 + e_3, -e_1 + e_2]$ of V corresponding to the elements $1, x - 2$ and $x^2 - 5x + 6$ in the diagonalized form of $xI - A$, respectively. The elements $f_1 = -7e_1 + 7e_2 + e_3$ and $f_2 = -e_1 + e_2$ are therefore $\mathbb{Q}[x]$ -module generators for the two cyclic factors of V in its invariant factor decomposition as a $\mathbb{Q}[x]$ -module. The corresponding \mathbb{Q} -vector space bases for these two factors are then f_1 and f_2 , $xf_2 = Tf_2$, i.e., $-7e_1 + 7e_2 + e_3$ and $-e_1 + e_2$, $T(-e_1 + e_2) = -4e_1 + 3e_2$. Then the matrix

$$P = \begin{pmatrix} -7 & -1 & -4 \\ 7 & 1 & 3 \\ 1 & 0 & 0 \end{pmatrix}$$

conjugates A into its rational canonical form:

$$P^{-1}AP = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 0 & -6 \\ 0 & 1 & 5 \end{pmatrix},$$

as one easily checks.

II. (Converting A Directly to Rational Canonical Form) We use the row operations involved in the diagonalization of $xI - A$ to determine the matrix P' of the algorithm above:

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \xrightarrow[\leftrightarrow C_2]{C_2-C_1} \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \xrightarrow{-C_1} \begin{pmatrix} -1 & -1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \rightarrow$$

$$C_1 - \xrightarrow[\leftrightarrow C_1]{(A-3I)C_2} \begin{pmatrix} 0 & -1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \xrightarrow[\leftrightarrow C_3]{C_3+7C_2} \begin{pmatrix} 0 & -1 & -7 \\ 0 & 1 & 7 \\ 0 & 0 & 1 \end{pmatrix} \xrightarrow{C_2 \leftrightarrow C_3} \begin{pmatrix} 0 & -7 & -1 \\ 0 & 7 & 1 \\ 0 & 1 & 0 \end{pmatrix} = P'.$$

Here we have $d_1 = 1$ and $d_2 = 2$, corresponding to the second and third nonzero columns of P' , respectively. The columns of P are therefore given by

$$\begin{pmatrix} -7 \\ 7 \\ 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, \quad A \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -4 \\ 3 \\ 0 \end{pmatrix},$$

respectively, which again gives the matrix P above.

- (3) For the 3×3 matrix A it was not necessary to perform the lengthy calculations above merely to determine the rational canonical form (equivalently, the invariant factors), as we saw in Example 1. For $n \times n$ matrices with $n \geq 4$, however, the computation of the characteristic and minimal polynomials is in general not sufficient for the determination of all the invariant factors, so the more extensive calculations of the previous example may become necessary. For example, consider the matrix

$$D = \begin{pmatrix} 1 & 2 & -4 & 4 \\ 2 & -1 & 4 & -8 \\ 1 & 0 & 1 & -2 \\ 0 & 1 & -2 & 3 \end{pmatrix}.$$

A short computation shows that the characteristic polynomial of D is $(x - 1)^4$. The possible minimal polynomials are then $x - 1$, $(x - 1)^2$, $(x - 1)^3$ and $(x - 1)^4$. Clearly $D - I \neq 0$ and another short computation shows that $(D - I)^2 = 0$, so the minimal polynomial for D is $(x - 1)^2$. There are then two possible sets of invariant factors:

$$x - 1, x - 1, (x - 1)^2 \quad \text{and} \quad (x - 1)^2, (x - 1)^2.$$

To determine the invariant factors for D we apply the procedure of the previous example to the 4×4 matrix

$$xI - D = \begin{pmatrix} x-1 & -2 & 4 & -4 \\ -2 & x+1 & -4 & 8 \\ -1 & 0 & x-1 & 2 \\ 0 & -1 & 2 & x-3 \end{pmatrix}.$$

The diagonal matrix obtained from this matrix by elementary row and column operations is the matrix

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & (x-1)^2 & 0 \\ 0 & 0 & 0 & (x-1)^2 \end{pmatrix},$$

which shows that the invariant factors for D are $(x - 1)^2, (x - 1)^2$ (one series of elementary row and column operations which diagonalize $xI - D$ are $R_1 \leftrightarrow R_3, -R_1$,

$$R_2 + 2R_1 \mapsto R_2, R_3 - (x-1)R_1 \mapsto R_3, C_3 + (x-1)C_1 \mapsto C_3, C_4 + 2C_1 \mapsto C_4, \\ R_2 \leftrightarrow R_4, -R_2, R_3 + 2R_2 \mapsto R_3, R_4 - (x+1)R_2 \mapsto R_4, C_3 + 2C_2 \mapsto C_3, \\ C_4 + (x-3)C_2 \mapsto C_4.$$

I. (Invariant Factor Decomposition) If e_1, e_2, e_3, e_4 is a basis for V in this case, then using the row operations in this diagonalization as in the previous example we see that the generators of V corresponding to the factors above are $(x-1)e_1 - 2e_2 - e_3 = 0$, $-2e_1 + (x+1)e_2 - e_4 = 0$, e_1, e_2 . Hence a vector space basis for the two direct factors in the invariant decomposition of V in this case is given by e_1, Te_1 and e_2, Te_2 where T is the linear transformation defined by D , i.e., $e_1, e_1 + 2e_2 + e_3$ and $e_2, 2e_1 - e_2 + e_4$. The corresponding matrix P relating these bases is

$$P = \begin{pmatrix} 1 & 1 & 0 & 2 \\ 0 & 2 & 1 & -1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

so that $P^{-1}DP$ is in rational canonical form:

$$P^{-1}DP = \begin{pmatrix} 0 & -1 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 2 \end{pmatrix}$$

as can easily be checked.

II. (Converting D Directly to Rational Canonical Form) As in Example 2 we determine the matrix P' of the algorithm from the row operations used in the diagonalization of $xI - D$:

$$\begin{array}{c} \left(\begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right) \xrightarrow[C_1 \leftrightarrow C_3]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right) \xrightarrow[-C_1]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right) \xrightarrow[]{} \\ \xrightarrow[C_1 - 2C_2]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ -2 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right) \xrightarrow[C_1 + (D-I)C_3]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{array} \right) \xrightarrow[C_2 \leftrightarrow C_4]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{array} \right) \xrightarrow[]{} \\ \xrightarrow[-C_2]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{array} \right) \xrightarrow[C_2 - 2C_3]{} \left(\begin{array}{cccc} 0 & -2 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{array} \right) \xrightarrow[C_2 + (D+I)C_4]{} \left(\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right) = P'. \end{array}$$

Here we have $d_1 = 2$ and $d_2 = 2$, corresponding to the third and fourth nonzero columns of P' . The columns of P are therefore given by

$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad D \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 1 \\ 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \quad D \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \\ 0 \\ 1 \end{pmatrix},$$

respectively, which again gives the matrix P above.

- (4) In this example we determine all similarity classes of matrices A with entries from \mathbb{Q} with characteristic polynomial $(x^4 - 1)(x^2 - 1)$. First note that any matrix with a degree

6 characteristic polynomial must be a 6×6 matrix. The polynomial $(x^4 - 1)(x^2 - 1)$ factors into irreducibles in $\mathbb{Q}[x]$ as $(x - 1)^2(x + 1)^2(x^2 + 1)$. Since the minimal polynomial $m_A(x)$ for A has the same roots as $c_A(x)$ it follows that $(x - 1)(x + 1)(x^2 + 1)$ divides $m_A(x)$. Suppose $a_1(x), \dots, a_m(x)$ are the invariant factors of some A , so $a_m(x) = m_A(x)$, $a_i(x) | a_{i+1}(x)$ (in particular, all the invariant factors divide $m_A(x)$) and $a_1(x)a_2(x) \cdots a_m(x) = (x^4 - 1)(x^2 - 1)$. One easily sees that the only permissible lists under these constraints are

- (a) $(x - 1)(x + 1), (x - 1)(x + 1)(x^2 + 1)$
- (b) $x - 1, (x - 1)(x + 1)^2(x^2 + 1)$
- (c) $x + 1, (x - 1)^2(x + 1)(x^2 + 1)$
- (d) $(x - 1)^2(x + 1)^2(x^2 + 1)$.

One can now easily write out the corresponding direct sums of companion matrices to obtain representatives of the 4 similarity classes. We shall see in the next section that there are still only 4 similarity classes even in $M_6(\mathbb{C})$.

- (5) In this example we find all similarity classes of 3×3 matrices A with entries from \mathbb{Q} satisfying $A^6 = I$. For each such A , its minimal polynomial divides $x^6 - 1$ and in $\mathbb{Q}[x]$ the complete factorization of this polynomial is

$$x^6 - 1 = (x - 1)(x + 1)(x^2 - x + 1)(x^2 + x + 1).$$

Conversely, if B is any 3×3 matrix whose minimal polynomial divides $x^6 - 1$, then $B^6 = I$. The only restriction on the minimal polynomial for B is that its degree is at most 3 (by the Cayley–Hamilton Theorem). The only possibilities for the minimal polynomial of such a matrix A are therefore

- (a) $x - 1$
- (b) $x + 1$
- (c) $x^2 - x + 1$
- (d) $x^2 + x + 1$
- (e) $(x - 1)(x + 1)$
- (f) $(x - 1)(x^2 - x + 1)$
- (g) $(x - 1)(x^2 + x + 1)$
- (h) $(x + 1)(x^2 - x + 1)$
- (i) $(x + 1)(x^2 + x + 1)$.

Under the constraints of the rational canonical form these give rise to the following permissible lists of invariant factors:

- (i) $x - 1, x - 1, x - 1$
- (ii) $x + 1, x + 1, x + 1$
- (iii) $x - 1, (x - 1)(x + 1)$
- (iv) $x + 1, (x - 1)(x + 1)$
- (v) $(x - 1)(x^2 - x + 1)$
- (vi) $(x - 1)(x^2 + x + 1)$
- (vii) $(x + 1)(x^2 - x + 1)$
- (viii) $(x + 1)(x^2 + x + 1)$.

Note that it is impossible to have a suitable set of invariant factors if the minimal polynomial is $x^2 + x + 1$ or $x^2 - x + 1$. One can now write out the corresponding

rational canonical forms; for example, (i) is I , (ii) is $-I$, and (iii) is

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

Note also that another way of phrasing this result is that any 3×3 matrix with entries from \mathbb{Q} whose order (multiplicatively, of course) divides 6 is similar to one of these 8 matrices, so this example determines all elements of orders 1, 2, 3 and 6 in the group $\mathrm{GL}_3(\mathbb{Q})$ (up to similarity).

EXERCISES

1. Prove that similar linear transformations of V (or $n \times n$ matrices) have the same characteristic and the same minimal polynomial.
2. Let M be as in Lemma 19. Prove that the minimal polynomial of M is the least common multiple of the minimal polynomials of A_1, \dots, A_k .
3. Prove that two 2×2 matrices over F which are not scalar matrices are similar if and only if they have the same characteristic polynomial.
4. Prove that two 3×3 matrices are similar if and only if they have the same characteristic and same minimal polynomials. Give an explicit counterexample to this assertion for 4×4 matrices.
5. Prove directly from the fact that the collection of *all* linear transformations of an n dimensional vector space V over F to itself form a vector space over F of dimension n^2 that the minimal polynomial of a linear transformation T has degree at most n^2 .
6. Prove that the constant term in the characteristic polynomial of the $n \times n$ matrix A is $(-1)^n \det A$ and that the coefficient of x^{n-1} is the negative of the sum of the diagonal entries of A (the sum of the diagonal entries of A is called the *trace* of A). Prove that $\det A$ is the product of the eigenvalues of A and that the trace of A is the sum of the eigenvalues of A .
7. Determine the eigenvalues of the matrix

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix}.$$

8. Verify that the characteristic polynomial of the companion matrix

$$\begin{pmatrix} 0 & 0 & 0 & \dots & 0 & -a_0 \\ 1 & 0 & 0 & \dots & 0 & -a_1 \\ 0 & 1 & 0 & \dots & 0 & -a_2 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & -a_{n-1} \end{pmatrix}$$

is

$$x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0.$$

9. Find the rational canonical forms of

$$\begin{pmatrix} 0 & -1 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{pmatrix}, \quad \begin{pmatrix} c & 0 & -1 \\ 0 & c & 1 \\ -1 & 1 & c \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 422 & 465 & 15 & -30 \\ -420 & -463 & -15 & 30 \\ 840 & 930 & 32 & -60 \\ -140 & -155 & -5 & 12 \end{pmatrix}.$$

10. Find all similarity classes of 6×6 matrices over \mathbb{Q} with minimal polynomial $(x+2)^2(x-1)$ (it suffices to give all lists of invariant factors and write out some of their corresponding matrices).
11. Find all similarity classes of 6×6 matrices over \mathbb{C} with characteristic polynomial $(x^4 - 1)(x^2 - 1)$.
12. Find all similarity classes of 3×3 matrices A over \mathbb{F}_2 satisfying $A^6 = I$ (compare with the answer we computed over \mathbb{Q}). Do the same for 4×4 matrices B satisfying $B^{20} = I$.
13. Prove that the number of similarity classes of 3×3 matrices over \mathbb{Q} with a given characteristic polynomial in $\mathbb{Q}[x]$ is the same as the number of similarity classes over any extension field of \mathbb{Q} . Give an example to show that this is not true in general for 4×4 matrices.
14. Determine all possible rational canonical forms for a linear transformation with characteristic polynomial $x^2(x^2 + 1)^2$.
15. Determine up to similarity all 2×2 rational matrices (i.e., $\in M_2(\mathbb{Q})$) of precise order 4 (multiplicatively, of course). Do the same if the matrix has entries from \mathbb{C} .
16. Show that $x^5 - 1 = (x - 1)(x^2 - 4x + 1)(x^2 + 5x + 1)$ in $\mathbb{F}_{19}[x]$. Use this to determine up to similarity all 2×2 matrices with entries from \mathbb{F}_{19} of (multiplicative) order 5.
17. Determine representatives for the conjugacy classes for $GL_3(\mathbb{F}_2)$. [Compare your answer with Theorem 15 and Proposition 14 of Chapter 6.]
18. Let V be a finite dimensional vector space over \mathbb{Q} and suppose T is a nonsingular linear transformation of V such that $T^{-1} = T^2 + T$. Prove that the dimension of V is divisible by 3. If the dimension of V is precisely 3 prove that all such transformations T are similar.
19. Let V be the infinite dimensional real vector space

$$\mathbb{R}^\infty = \{(a_0, a_1, a_2, \dots) \mid a_0, a_1, a_2, \dots \in \mathbb{R}\}.$$

Define the map $T : V \rightarrow V$ by $T(a_0, a_1, a_2, \dots) = (0, a_0, a_1, a_2, \dots)$. Prove that T has no eigenvectors.

20. Let ℓ be a prime and let $\Phi_\ell(x) = \frac{x^\ell - 1}{x - 1} = x^{\ell-1} + x^{\ell-2} + \dots + x + 1 \in \mathbb{Z}[x]$ be the ℓ^{th} cyclotomic polynomial, which is irreducible over \mathbb{Q} (Example 4 following Corollary 9.14). This exercise determines the smallest degree of a factor of $\Phi_\ell(x)$ modulo p for any prime p and so in particular determines when $\Phi_\ell(x)$ is irreducible modulo p . (This actually determines the complete factorization of $\Phi_\ell(x)$ modulo p — cf. Exercise 8 of Section 13.6.)
- (a) Show that if $p = \ell$ then $\Phi_\ell(x)$ is divisible by $x - 1$ in $\mathbb{F}_\ell[x]$.
- (b) Suppose $p \neq \ell$ and let f denote the order of p in \mathbb{F}_ℓ^\times , i.e., f is the smallest power of p with $p^f \equiv 1 \pmod{\ell}$. Show that $m = f$ is the first value of m for which the group $GL_m(\mathbb{F}_p)$ contains an element A of order ℓ . [Use the formula for the order of this group at the end of Section 11.1.]
- (c) Show that $\Phi_\ell(x)$ is not divisible by any polynomial of degree smaller than f in $\mathbb{F}_p[x]$ [consider the companion matrix for such a divisor and use (b)]. Let $m_A(x) \in \mathbb{F}_p[x]$ denote the minimal polynomial for the matrix A in (b) and conclude that $m_A(x)$ is irreducible of degree f and divides $\Phi_\ell(x)$ in $\mathbb{F}_p[x]$.

- (d) In particular, prove that $\Phi_\ell(x)$ is irreducible modulo p if and only if $\ell - 1$ is the smallest power of p which is congruent to 1 modulo ℓ , i.e., p is a primitive root modulo ℓ .
21. Prove that the first two elementary row and column operations described before Theorem 21 do not change the determinant of the matrix and the third elementary operation multiplies the determinant by a unit. Conclude from Theorem 21 that the characteristic polynomial of A differs by a unit from the product of the invariant factors of A . Since both these polynomials are monic by definition, conclude that they are equal (this gives an alternate proof of Proposition 20).

The following exercises outline the proof of Theorem 21. They carry out explicitly the construction described in Exercises 16 to 19 of the previous section for the Euclidean Domain $F[x]$. Let V be an n -dimensional vector space with basis v_1, v_2, \dots, v_n and let T be the linear transformation of V defined by the matrix A and this choice of basis, i.e., T is the linear transformation with

$$T(v_j) = \sum_{i=1}^n a_{ij} v_i, \quad j = 1, 2, \dots, n$$

where $A = (a_{ij})$. Let $F[x]^n$ be the free module of rank n over $F[x]$ and let $\xi_1, \xi_2, \dots, \xi_n$ denote a basis. Then we have a natural surjective $F[x]$ -module homomorphism

$$\varphi : F[x]^n \rightarrow V$$

defined by mapping ξ_i to v_i , $i = 1, 2, \dots, n$. As indicated in the exercises of the previous section the invariant factors for the $F[x]$ -module V can be determined once we have determined a set of generators and the corresponding relations matrix for $\ker \varphi$. Since by definition x acts on V by the linear transformation T , we have

$$x(v_j) = \sum_{i=1}^n a_{ij} v_i, \quad j = 1, 2, \dots, n.$$

22. Show that the elements

$$v_j = -a_{1j}\xi_1 - \cdots - a_{j-1,j}\xi_{j-1} + (x - a_{jj})\xi_j - a_{j+1,j}\xi_{j+1} - \cdots - a_{nj}\xi_n$$

for $j = 1, 2, \dots, n$ are elements of the kernel of φ .

23. (a) Show that $x\xi_j = v_j + f_j$ where $f_j \in F\xi_1 + \cdots + F\xi_n$ is an element in the F -vector space spanned by ξ_1, \dots, ξ_n .
(b) Show that

$$F[x]\xi_1 + \cdots + F[x]\xi_n = (F[x]v_1 + \cdots + F[x]v_n) + (F\xi_1 + \cdots + F\xi_n).$$

24. Show that v_1, v_2, \dots, v_n generate the kernel of φ . [Use the previous result to show that any element of $\ker \varphi$ is the sum of an element in the module generated by v_1, v_2, \dots, v_n and an element of the form $b_1\xi_1 + \cdots + b_n\xi_n$ where the b_i are elements of F . Then show that such an element is in $\ker \varphi$ if and only if all the b_i are 0 since v_1, \dots, v_n are a basis for V over F .]

25. Show that the generators v_1, v_2, \dots, v_n of $\ker \varphi$ have corresponding relations matrix

$$\begin{pmatrix} x - a_{11} & -a_{21} & \cdots & -a_{n1} \\ -a_{12} & x - a_{22} & \cdots & -a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ -a_{1n} & -a_{2n} & \cdots & x - a_{nn} \end{pmatrix} = xI - A^t,$$

where A' is the transpose of A . Conclude that Theorem 21 and the algorithm for determining the invariant factors of A follows by Exercises 16 to 19 in the previous section (note that the row and column operations necessary to diagonalize this relations matrix are the column and row operations necessary to diagonalize the matrix in Theorem 21, which explains why the invariant factor algorithm keeps track of the *row* operations used).

12.3 THE JORDAN CANONICAL FORM

We continue with the notation in the previous section: F is a field, $F[x]$ is the ring of polynomials in x with coefficients in F , V is a finite dimensional vector space over F of dimension n , T is a fixed linear transformation of V by which we make V into an $F[x]$ -module, and A is an $n \times n$ matrix with coefficients in F . Recall that once a basis for V has been fixed any linear transformation T defines a matrix A and conversely any matrix A defines a linear transformation T .

In the previous section we used the invariant factor form of the Fundamental Theorem for finitely generated modules over the Principal Ideal Domain $F[x]$ to obtain the rational canonical form for such a linear transformation T and the rational canonical form for such an $n \times n$ matrix A . In this section we use the elementary divisor form of the Fundamental Theorem to obtain the *Jordan canonical form*. We shall see that matrices in this canonical form are as close to being diagonal matrices as possible, so the matrices are simpler than in the rational canonical form (but we lose some of the “rationality” results).

The elementary divisors of a module are the prime power divisors of its invariant factors (this was Corollary 10). For the $F[x]$ -module V the invariant factors were monic polynomials $a_1(x), a_2(x), \dots, a_m(x)$ of degree at least one (with $a_1(x) | a_2(x) | \dots | a_m(x)$), so the associated elementary divisors are the powers of the irreducible polynomial factors of these polynomials. These polynomials are only defined up to multiplication by a unit and, as in the case of the invariant factors, we can specify them uniquely by requiring that they be monic.

To obtain the simplest possible elementary divisors we shall assume that the polynomials $a_1(x), a_2(x), \dots, a_m(x)$ factor completely into linear factors, i.e., that the elementary divisors of V are powers $(x - \lambda)^k$ of linear polynomials. Since the product of the elementary divisors is the characteristic polynomial, this is equivalent to the assumption that the field F contains all the eigenvalues of the linear transformation T (equivalently, of the matrix A representing the linear transformation T).

Under this assumption on F , it follows immediately from Theorem 6 that V is the direct sum of finitely many cyclic $F[x]$ -modules of the form $F[x]/(x - \lambda)^k$ where $\lambda \in F$ is one of the eigenvalues of T , corresponding to the elementary divisors of V .

We now choose a vector space basis for each of the direct summands corresponding to the elementary divisors of V for which the corresponding matrix for T is particularly simple. Recall that by definition of the $F[x]$ -module structure the linear transformation T acting on V is the element x acting by multiplication on each of the direct summands $F[x]/(x - \lambda)^k$.

Consider the elements

$$(\bar{x} - \lambda)^{k-1}, (\bar{x} - \lambda)^{k-2}, \dots, \bar{x} - \lambda, 1,$$

in the quotient $F[x]/(x - \lambda)^k$. Expanding each of these polynomials in \bar{x} we see that the matrix relating these elements to the F -basis $\bar{x}^{k-1}, \bar{x}^{k-2}, \dots, \bar{x}, 1$ of $F[x]/(x - \lambda)^k$ is upper triangular with 1's along the diagonal. Since this is an invertible matrix (having determinant 1), it follows that the elements above are an F -basis for $F[x]/(x - \lambda)^k$. With respect to this basis the linear transformation of multiplication by x acts in a particularly simple manner (note that $x = \lambda + (x - \lambda)$ and that $(\bar{x} - \lambda)^k = 0$ in the quotient):

$$\begin{array}{ll} (\bar{x} - \lambda)^{k-1} & \mapsto \lambda \cdot (\bar{x} - \lambda)^{k-1} + (\bar{x} - \lambda)^k = \lambda \cdot (\bar{x} - \lambda)^{k-1} \\ (\bar{x} - \lambda)^{k-2} & \mapsto \lambda \cdot (\bar{x} - \lambda)^{k-2} + (\bar{x} - \lambda)^{k-1} \\ x : & \vdots \\ \bar{x} - \lambda & \mapsto \lambda \cdot (\bar{x} - \lambda) + (\bar{x} - \lambda)^2 \\ 1 & \mapsto \lambda \cdot 1 + (\bar{x} - \lambda). \end{array}$$

With respect to this basis, the matrix for multiplication by x is therefore

$$\begin{pmatrix} \lambda & 1 & & & \\ & \lambda & \ddots & & \\ & & \ddots & 1 & \\ & & & \lambda & 1 \\ & & & & \lambda \end{pmatrix}$$

where the blank entries are all zero. Such matrices are given a name:

Definition. The $k \times k$ matrix with λ along the main diagonal and 1 along the first superdiagonal depicted above is called the $k \times k$ *elementary Jordan matrix with eigenvalue λ* or the *Jordan block of size k with eigenvalue λ* .

Applying this to each of the cyclic factors of V in its elementary divisor decomposition we obtain a vector space basis for V with respect to which the linear transformation T has as matrix the direct sum of the Jordan blocks corresponding to the elementary divisors of V , i.e., is block diagonal with Jordan blocks along the diagonal:

$$\begin{pmatrix} J_1 & & & \\ & J_2 & & \\ & & \ddots & \\ & & & J_t \end{pmatrix}.$$

Notice that this matrix is uniquely determined up to permutation of the blocks along the diagonal by the elementary divisors of the $F[x]$ -module V and conversely, by Theorem 9, the list of elementary divisors uniquely determines the module V up to $F[x]$ -module isomorphism.

Definition.

- (1) A matrix is said to be in *Jordan canonical form* if it is a block diagonal matrix with Jordan blocks along the diagonal.
- (2) A *Jordan canonical form* for a linear transformation T is a matrix representing T which is in Jordan canonical form.

We have proved that any linear transformation T has a Jordan canonical form. As in the case of the rational canonical form, it follows from the uniqueness of the elementary divisors that the Jordan canonical form is unique up to a permutation of the Jordan blocks along the diagonal (hence is called *the* Jordan canonical form for T). We summarize this in the following theorem.

Theorem 22. (Jordan Canonical Form for Linear Transformations) Let V be a finite dimensional vector space over the field F and let T be a linear transformation of V . Assume F contains all the eigenvalues of T .

- (1) There is a basis for V with respect to which the matrix for T is in Jordan canonical form, i.e., is a block diagonal matrix whose diagonal blocks are the Jordan blocks for the elementary divisors of V .
- (2) The Jordan canonical form for T is unique up to a permutation of the Jordan blocks along the diagonal.

As for the rational canonical form, the following theorem gives the corresponding statement for $n \times n$ matrices over F .

Theorem 23. (Jordan Canonical Form for Matrices) Let A be an $n \times n$ matrix over the field F and assume F contains all the eigenvalues of A .

- (1) The matrix A is similar to a matrix in Jordan canonical form, i.e., there is an invertible $n \times n$ matrix P over F such that $P^{-1}AP$ is a block diagonal matrix whose diagonal blocks are the Jordan blocks for the elementary divisors of A .
- (2) The Jordan canonical form for A is unique up to a permutation of the Jordan blocks along the diagonal.

The Jordan canonical form differs from a diagonal matrix only by the possible presence of some 1's along the first superdiagonal (and then only if there are Jordan blocks of size greater than one), hence is close to being a diagonal matrix. The following result shows in particular that the Jordan canonical form for a matrix A is as close to being a diagonal matrix as possible.

Corollary 24.

- (1) If a matrix A is similar to a diagonal matrix D , then D is the Jordan canonical form of A .
- (2) Two diagonal matrices are similar if and only if their diagonal entries are the same up to a permutation.

Proof: The first assertion is immediate from the uniqueness of Jordan canonical forms because a diagonal matrix is itself in Jordan form (with Jordan blocks of size 1). The uniqueness of the Jordan canonical form gives (2).

The next corollary gives a criterion to determine when a matrix A can be diagonalized.

Corollary 25. If A is an $n \times n$ matrix with entries from F and F contains all the eigenvalues of A , then A is similar to a diagonal matrix over F if and only if the minimal polynomial of A has no repeated roots.

Proof: Suppose A is similar to a diagonal matrix. The minimal polynomial of a diagonal matrix has no repeated roots (its roots are precisely the distinct elements along the diagonal). Since similar matrices have the same minimal polynomial it follows that the minimal polynomial for A has no repeated roots.

Conversely, suppose the minimal polynomial for A has no repeated roots and let B be the Jordan canonical form of A . The matrix B is a block diagonal matrix with elementary Jordan matrices down the diagonal. By the exercises at the end of the preceding section the minimal polynomial for B is the least common multiple of the minimal polynomials of the Jordan blocks. It is easy to see directly that a Jordan block of size k with eigenvalue λ has minimal polynomial $(x - \lambda)^k$ (note that this is immediate from the fact that each elementary Jordan matrix gives the action on a *cyclic* $F[x]$ -submodule whose annihilator is $(x - \lambda)^k$). Since A and B have the same minimal polynomial, the least common multiple of the $(x - \lambda)^k$ cannot have any repeated roots. It follows that k must be 1, i.e., that each Jordan block must be of size one and B is a diagonal matrix.

Changing From One Canonical Form to Another

We continue to assume that the field F contains all the eigenvalues of T (or A) so both the rational and Jordan canonical forms exist over F . The process of passing from one form to the other is exactly the same algorithm described in Section 5.2 for finite abelian groups (where the elementary divisors were determined from the list of invariant factors and vice versa).

In brief summary, recall that the elementary divisors are the prime power divisors of the invariant factors. They are obtained from the invariant factors by writing each invariant factor as a product of distinct linear factors to powers; the resulting set of powers of linear polynomials is the set of elementary divisors. For example, if the invariant factors of T are

$$(x - 1)(x - 3)^3, \quad (x - 1)(x - 2)(x - 3)^3, \quad (x - 1)(x - 2)^2(x - 3)^3$$

then the elementary divisors are

$$(x - 1), \quad (x - 3)^3, \quad (x - 1), \quad (x - 2), \quad (x - 3)^3, \quad (x - 1), \quad (x - 2)^2, \quad (x - 3)^3.$$

The largest invariant factor is the product of the largest of the distinct prime powers among the elementary divisors, the next largest invariant factor is the product of the largest of the distinct prime powers among the remaining elementary divisors, and so on. Given a list of elementary divisors we can find the list of invariant factors by first arranging the elementary divisors into n separate lists, one for each eigenvalue. In each of these n lists arrange the polynomials in increasing (i.e., nondecreasing) degree. Next arrange for all n lists to have the same length by appending an appropriate number of the constant polynomial 1. Now form the i^{th} invariant factor by taking the product of

the i^{th} polynomial in each of these lists. For example, if the elementary divisors of T are

$$(x - 1)^3, (x + 4), (x + 4)^2, (x - 5)^2, (x - 1)^5, (x - 1)^3, (x - 5)^3, (x - 1)^4, (x + 4)^3$$

then the intermediate lists are

$$\begin{array}{llll} (1) & (x - 1)^3, & (x - 1)^3, & (x - 1)^4, & (x - 1)^5 \\ (2) & 1, & x + 4, & (x + 4)^2, & (x + 4)^3 \\ (3) & 1, & 1, & (x - 5)^2, & (x - 5)^3 \end{array}$$

so the list of invariant factors is

$$(x - 1)^3, (x - 1)^3(x + 4), (x - 1)^4(x + 4)^2(x - 5)^2, (x - 1)^5(x + 4)^3(x - 5)^3.$$

Elementary Divisor Decomposition Algorithm: Converting to Jordan Canonical Forms

Theorem 21 indicates a computational procedure to determine the invariant factors of any given matrix A . Factorization of these invariant factors produces the elementary divisors of A , hence determines the Jordan canonical form for A as above.

The Invariant Factor Decomposition Algorithm following Theorem 21 starts with a basis e_1, \dots, e_n for V and produces a set f_1, \dots, f_m of elements of V which are $F[x]$ -module generators for the cyclic factors in the invariant factor decomposition of V (with annihilators $(a_1(x)), \dots, (a_m(x))$, respectively). Since the elementary divisor decomposition is obtained from the invariant factor decomposition by applying the Chinese Remainder Theorem to the cyclic modules $F[x]/(a_i(x))$, this gives a set of $F[x]$ -module generators for the cyclic factors in the elementary divisor decomposition of V . These elements then give rise to an explicit vector space basis for V with respect to which the linear transformation corresponding to A is in Jordan canonical form (equivalently, an explicit matrix P such that $P^{-1}AP$ is in Jordan canonical form). As for the Invariant Factor Decomposition Algorithm we state the result first in the general context of decomposing a vector space and then describe the algorithm to convert a given $n \times n$ matrix A to Jordan canonical form.

Explicit numerical examples of this algorithm are given later in Examples 2 and 3.

Elementary Divisor Decomposition Algorithm

(1) to (3): The first three steps in the algorithm are those from the Invariant Factor Decomposition Algorithm following Theorem 21.

(4) For each invariant factor $a(x)$ computed for A write

$$a(x) = (x - \lambda_1)^{\alpha_1}(x - \lambda_2)^{\alpha_2} \dots (x - \lambda_s)^{\alpha_s}$$

where $\lambda_1, \dots, \lambda_s \in F$ are distinct. Let $f \in V$ be the $F[x]$ -module generator for the cyclic factor corresponding to the invariant factor $a(x)$ computed in (3). Then the elements

$$\frac{a(x)}{(x - \lambda_1)^{\alpha_1}}f, \quad \frac{a(x)}{(x - \lambda_2)^{\alpha_2}}f, \quad \dots, \quad \frac{a(x)}{(x - \lambda_s)^{\alpha_s}}f$$

(note that the $\frac{a(x)}{(x - \lambda_i)^{\alpha_i}} \in F[x]$ are polynomials) are $F[x]$ -module generators for the cyclic factors of V corresponding to the elementary divisors

$$(x - \lambda_1)^{\alpha_1}, \quad (x - \lambda_2)^{\alpha_2}, \quad \dots, \quad (x - \lambda_s)^{\alpha_s},$$

respectively.

- (5) If $g_i = \frac{a(x)}{(x - \lambda_i)^{\alpha_i}} f$ is the $F[x]$ -module generator for the cyclic factor of V corresponding to the elementary divisor $(x - \lambda_i)^{\alpha_i}$ then the corresponding *vector space* basis for this cyclic factor of V is given by the elements

$$(T - \lambda_i)^{\alpha_i-1} g_i, \quad (T - \lambda_i)^{\alpha_i-2} g_i, \quad \dots, \quad (T - \lambda_i) g_i, \quad g_i.$$

- (6) Write the k^{th} element of the vector space basis computed in (5) in terms of the original vector space basis $[e_1, e_2, \dots, e_n]$ for V and use the coordinates for the k^{th} column of an $n \times n$ matrix P . Then $P^{-1}AP$ is in Jordan canonical form (with Jordan blocks appearing in the order used in (5) for the cyclic factors of V).

Converting an $n \times n$ Matrix to Jordan Canonical Form

- (1) to (2): The first two steps are those from the algorithm for Converting an $n \times n$ matrix to Rational Canonical Form following Theorem 21.
 (3) When $xI - A$ has been diagonalized to the form in Theorem 21 the first $n-m$ columns of the matrix P' are 0 (providing a useful numerical check on the computations) and the remaining m columns of P' are nonzero. For each successive $i = 1, 2, \dots, m$:

- (a) Factor the i^{th} nonconstant diagonal element (which is of degree d_i):

$$a(x) = (x - \lambda_1)^{\alpha_1}(x - \lambda_2)^{\alpha_2} \dots (x - \lambda_s)^{\alpha_s}$$

where $\lambda_1, \dots, \lambda_s \in F$ are distinct (here $a(x) = a_i(x)$ is the i^{th} nonconstant diagonal element and s depends on i).

- (b) Multiply the i^{th} nonzero column of P' successively by the d_i matrices:

$$\begin{array}{ccccccc} (A - \lambda_1 I)^{\alpha_1-1} & (A - \lambda_2 I)^{\alpha_2} & \dots & (A - \lambda_s I)^{\alpha_s} \\ (A - \lambda_1 I)^{\alpha_1-2} & (A - \lambda_2 I)^{\alpha_2} & \dots & (A - \lambda_s I)^{\alpha_s} \\ \vdots & & & & & & \\ (A - \lambda_1 I)^0 & (A - \lambda_2 I)^{\alpha_2} & \dots & (A - \lambda_s I)^{\alpha_s} \\ \\ (A - \lambda_1 I)^{\alpha_1} & (A - \lambda_2 I)^{\alpha_2-1} & \dots & (A - \lambda_s I)^{\alpha_s} \\ (A - \lambda_1 I)^{\alpha_1} & (A - \lambda_2 I)^{\alpha_2-2} & \dots & (A - \lambda_s I)^{\alpha_s} \\ \vdots & & & & & & \\ (A - \lambda_1 I)^{\alpha_1} & (A - \lambda_2 I)^0 & \dots & (A - \lambda_s I)^{\alpha_s} \\ \vdots & & & & & & \end{array}$$

$$\begin{array}{c}
 \vdots \\
 (A - \lambda_1 I)^{\alpha_1} (A - \lambda_2 I)^{\alpha_2} \dots (A - \lambda_s I)^{\alpha_s - 1} \\
 (A - \lambda_1 I)^{\alpha_1} (A - \lambda_2 I)^{\alpha_2} \dots (A - \lambda_s I)^{\alpha_s - 2} \\
 \vdots \\
 (A - \lambda_1 I)^{\alpha_1} (A - \lambda_2 I)^{\alpha_2} \dots (A - \lambda_s I)^0.
 \end{array}$$

- (c) Use the column vectors resulting from (b) (in that order) as the next d_i columns of an $n \times n$ matrix P .

Then $P^{-1}AP$ is in Jordan canonical form (whose Jordan blocks correspond to the ordering of the factors in (a)).

Examples

We can use Jordan canonical forms to carry out the same analysis of matrices that we did as examples of the use of rational canonical forms. In some instances, when the field is enlarged, the number of similarity classes increases (the number of similarity classes can never decrease when we extend the field by Corollary 18(2)).

- (1) Let A , B and C be the matrices in Example 1 of the previous section and let $F = \mathbb{Q}$. Note that \mathbb{Q} contains all the eigenvalues for these matrices. Since we have already determined the invariant factors of these matrices we can immediately obtain their elementary divisors. The elementary divisors of A are $x - 2$, $x - 2$ and $x - 3$ and the elementary divisors of B and C are $(x - 2)^2$ and $x - 3$ so the respective Jordan canonical forms are:

$$\begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix}.$$

Notice that A is similar to a diagonal matrix but, by Corollary 25, B and C are not.

- (2) For the matrix A , we determined in Example 2 of the previous section that $f_1 = -7e_1 + 7e_2 + e_3$ and $f_2 = -e_1 + e_2$ were $\mathbb{Q}[x]$ -module generators for the two cyclic factors of V in its invariant factor decomposition, corresponding to the invariant factors $x - 2$ and $(x - 2)(x - 3)$, respectively. Using the first algorithm described above, the elements f_1 , $(x - 3)f_2$ and $(x - 2)f_2$ are therefore $\mathbb{Q}[x]$ -module generators for the three cyclic factors of V in its elementary divisor decomposition, corresponding to the elementary divisors $x - 2$, $x - 2$, and $x - 3$. An easy computation shows that these are the elements $-7e_1 + 7e_2 + e_3$, $-e_1$ and $-2e_1 + e_2$, respectively. Then the matrix

$$P = \begin{pmatrix} -7 & -1 & -2 \\ 7 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

conjugates A into its Jordan canonical form:

$$P^{-1}AP = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix},$$

as one easily checks.

The columns of this matrix can also be obtained following the second algorithm above, using the nonzero columns of the matrix P' computed in Example 2 of the

previous section:

$$(A - 2I)^0 \begin{pmatrix} -7 \\ 7 \\ 1 \end{pmatrix} = \begin{pmatrix} -7 \\ 7 \\ 1 \end{pmatrix}$$

and

$$(A - 2I)^0(A - 3I)^1 \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ 0 \\ 0 \end{pmatrix}, \quad (A - 2I)^1(A - 3I)^0 \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -2 \\ 1 \\ 0 \end{pmatrix},$$

respectively, which again gives the matrix P .

- (3) For the 4×4 matrix D of Example 3 of the previous section, the invariant factors were $(x - 1)^2, (x - 1)^2$, with corresponding $\mathbb{Q}[x]$ -module generators $f_1 = e_1$ and $f_2 = e_2$, respectively. These are also the elementary divisors for this matrix. The corresponding vector space bases for these two factors are given by $(T - 1)f_1, f_1$ and $(T - 1)f_2, f_2$, respectively. An easy computation shows these are the elements $2e_2 + e_3, e_1$ and $2e_1 - e_2 + e_4, e_2$, respectively. Then the matrix

$$P = \begin{pmatrix} 0 & 1 & 2 & 0 \\ 2 & 0 & -2 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

conjugates D into its Jordan canonical form:

$$P^{-1}DP = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

as can easily be checked.

The columns of this matrix can also be obtained following the second algorithm above, using the nonzero columns of the matrix P' computed in Example 3 of the previous section:

$$(D - I)^1 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 2 \\ 1 \\ 0 \end{pmatrix}, \quad (D - I)^0 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix},$$

and

$$(D - I)^1 \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 2 \\ -2 \\ 0 \\ 1 \end{pmatrix}, \quad (D - I)^0 \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix},$$

respectively, which again gives the matrix P .

- (4) The set of similarity classes of 6×6 matrices with entries from \mathbb{C} with characteristic polynomial $(x^4 - 1)(x^2 - 1)$ consists of the 4 classes represented by the rational canonical forms in the preceding set of examples (there are no additional lists of invariant factors over \mathbb{C}). Their Jordan canonical forms cannot all be written over \mathbb{Q} , however. For instance, if the invariant factors are

$$(x - 1)(x + 1) \quad \text{and} \quad (x - 1)(x + 1)(x^2 + 1)$$

then the elementary divisors are

$$x - 1, \quad x + 1, \quad x - 1, \quad x + 1, \quad x - i, \quad x + i,$$

where i is a square root of -1 in \mathbb{C} , so the Jordan form for this matrix is a diagonal matrix with diagonal entries $1, 1, -1, -1, i, -i$.

- (5) In contrast, the set of similarity classes of 3×3 matrices, A , over \mathbb{C} satisfying $A^6 = I$ is considerably larger than that over \mathbb{Q} . If A is any such matrix, $m_A(x) \mid x^6 - 1$ so since the latter polynomial has no repeated roots in \mathbb{C} , the minimal polynomial of A has no repeated roots. By Corollary 25 the Jordan canonical form of A is a diagonal matrix. Since this diagonal matrix has the same minimal polynomial, its 6th power is also the identity, and so each diagonal entry is a 6th root of unity. For each list $\zeta_1, \zeta_2, \zeta_3$ of 6th roots of unity we obtain a Jordan canonical form, and two such forms are the same (i.e., give rise to similar matrices) if and only if the lists are permuted versions of each other. One finds that there are, up to similarity, 56 classes of such A 's.

EXERCISES

- Suppose the vector space V is the direct sum of cyclic $F[x]$ -modules whose annihilators are $(x + 1)^2$, $(x - 1)(x^2 + 1)^2$, $(x^4 - 1)$ and $(x + 1)(x^2 - 1)$. Determine the invariant factors and elementary divisors for V .
- Prove that if $\lambda_1, \dots, \lambda_n$ are the eigenvalues of the $n \times n$ matrix A then $\lambda_1^k, \dots, \lambda_n^k$ are the eigenvalues of A^k for any $k \geq 0$.
- Use the method of Example 2 above to determine explicit matrices P_1 and P_2 with $P_1^{-1}BP_1$ and $P_2^{-1}CP_2$ in Jordan canonical form. Use this to explicitly construct a matrix Q which conjugates B into C (proving directly that these matrices are similar).
- Prove that the Jordan canonical form for the matrix

$$\begin{pmatrix} 9 & 4 & 5 \\ -4 & 0 & -3 \\ -6 & -4 & -2 \end{pmatrix}$$

is that stated at the beginning of this chapter. Explicitly determine a matrix P which conjugates this matrix to its Jordan canonical form. Explain why this matrix cannot be diagonalized.

- Compute the Jordan canonical form for the matrix

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & -2 \\ 0 & 1 & 3 \end{pmatrix}.$$

- Determine which of the following matrices are similar:

$$\begin{pmatrix} -1 & 4 & -4 \\ 2 & -1 & 3 \\ 0 & -4 & 3 \end{pmatrix} \quad \begin{pmatrix} -3 & -4 & 0 \\ 2 & 3 & 0 \\ 8 & 8 & 1 \end{pmatrix} \quad \begin{pmatrix} -3 & 2 & -4 \\ 2 & 1 & 0 \\ 3 & -1 & 3 \end{pmatrix} \quad \begin{pmatrix} -1 & 4 & -4 \\ 0 & -3 & 2 \\ 0 & -4 & 3 \end{pmatrix}.$$

- Determine the Jordan canonical forms for the following matrices:

$$\begin{pmatrix} 5 & 4 & 1 \\ -1 & 0 & 0 \\ -3 & -4 & 1 \end{pmatrix} \quad \begin{pmatrix} 3 & 4 & 2 \\ -2 & -3 & -1 \\ -4 & -4 & -3 \end{pmatrix}.$$

8. Prove that the matrices

$$A = \begin{pmatrix} 5 & 6 & 0 \\ -3 & -4 & 0 \\ -2 & 0 & 1 \end{pmatrix} \quad B = \begin{pmatrix} 3 & -1 & 2 \\ -10 & 6 & -14 \\ -6 & 3 & -7 \end{pmatrix}$$

are similar. Prove that both A and B can be diagonalized and determine explicit matrices P_1 and P_2 with $P_1^{-1}AP_1$ and $P_2^{-1}BP_2$ in diagonal form.

9. Prove that the matrices

$$A = \begin{pmatrix} -8 & -10 & -1 \\ 7 & 9 & 1 \\ 3 & 2 & 0 \end{pmatrix} \quad B = \begin{pmatrix} -3 & 2 & -4 \\ 4 & -1 & 4 \\ 4 & -2 & 5 \end{pmatrix}$$

both have $(x - 1)^2(x + 1)$ as characteristic polynomial but that one can be diagonalized and the other cannot. Determine the Jordan canonical form for both matrices.

10. Find all Jordan canonical forms of 2×2 , 3×3 and 4×4 matrices over \mathbb{C} .

11. Verify that the characteristic polynomial of

$$A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ -2 & -2 & 0 & 1 \\ -2 & 0 & -1 & -2 \end{pmatrix}$$

is a product of linear factors over \mathbb{Q} . Determine the rational and Jordan canonical forms for A over \mathbb{Q} .

12. Determine the Jordan canonical form for the matrix

$$\begin{pmatrix} 1 & 2 & 0 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

13. Determine the Jordan canonical form for the matrix

$$\begin{pmatrix} 3 & 0 & -2 & -3 \\ 4 & -8 & 14 & -15 \\ 2 & -4 & 7 & -7 \\ 0 & 2 & -4 & 3 \end{pmatrix}.$$

14. Prove that the matrices

$$A = \begin{pmatrix} 2 & 0 & 0 & 0 \\ -4 & -1 & -4 & 0 \\ 2 & 1 & 3 & 0 \\ -2 & 4 & 9 & 1 \end{pmatrix} \quad B = \begin{pmatrix} 5 & 0 & -4 & -7 \\ 3 & -8 & 15 & -13 \\ 2 & -4 & 7 & -7 \\ 1 & 2 & -5 & 1 \end{pmatrix}$$

are similar.

15. Prove that the matrices

$$A = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \quad B = \begin{pmatrix} 5 & 2 & -8 & -8 \\ -6 & -3 & 8 & 8 \\ -3 & -1 & 3 & 4 \\ 3 & 1 & -4 & -5 \end{pmatrix}$$

both have characteristic polynomial $(x - 3)(x + 1)^3$. Determine whether they are similar and determine the Jordan canonical form for each matrix.

- 16.** Determine the Jordan canonical form for the matrix

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

and determine a matrix P which conjugates this matrix into its Jordan canonical form.

- 17.** Prove that any matrix A is similar to its transpose A' .
- 18.** Determine all possible Jordan canonical forms for a linear transformation with characteristic polynomial $(x - 2)^3(x - 3)^2$.
- 19.** Prove that all $n \times n$ matrices with characteristic polynomial $f(x)$ are similar if and only if $f(x)$ has no repeated factors in its unique factorization in $F[x]$.
- 20.** Show that the following matrices are similar in $M_p(\mathbb{F}_p)$ ($p \times p$ matrices with entries from \mathbb{F}_p):

$$\begin{pmatrix} 0 & 0 & 0 & \dots & 0 & 1 \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 & 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & 1 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 1 \\ 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix}.$$

- 21.** Show that if $A^2 = A$ then A is similar to a diagonal matrix which has only 0's and 1's along the diagonal.
- 22.** Prove that an $n \times n$ matrix A with entries from \mathbb{C} satisfying $A^3 = A$ can be diagonalized. Is the same statement true over *any* field F ?
- 23.** Suppose A is a 2×2 matrix with entries from \mathbb{Q} for which $A^3 = I$ but $A \neq I$. Write A in rational canonical form and in Jordan canonical form viewed as a matrix over \mathbb{C} .
- 24.** Prove there are no 3×3 matrices A over \mathbb{Q} with $A^8 = I$ but $A^4 \neq I$.
- 25.** Determine the Jordan canonical form for the $n \times n$ matrix over \mathbb{Q} whose entries are all equal to 1.
- 26.** Determine the Jordan canonical form for the $n \times n$ matrix over \mathbb{F}_p whose entries are all equal to 1 (the answer depends on whether or not p divides n).
- 27.** Determine the Jordan canonical form for the $n \times n$ matrix over \mathbb{Q} whose entries are all equal to 1 except that the entries along the main diagonal are all equal to 0.
- 28.** Determine the Jordan canonical form for the $n \times n$ matrix over \mathbb{F}_p whose entries are all equal to 1 except that the entries along the main diagonal are all equal to 0.

The direct sum of the cyclic submodules of V corresponding to all the elementary divisors of V which are powers of the same $x - \lambda$ is called the *generalized eigenspace* of T corresponding to the eigenvalue λ . Note that this is the p -primary component of V for the prime $p = x - \lambda$ of $F[x]$ and consists of the elements of V which are annihilated by some power of the linear transformation $T - \lambda$. The matrix for T on the generalized eigenspace for λ is the block diagonal matrix of all Jordan blocks for T with the same eigenvalue λ .

- 29.** Suppose V_i is the generalized eigenspace of T corresponding to eigenvalue λ_i . For any $k \geq 0$, prove that the nullity of $T - \lambda_i$ on the subspace $(T - \lambda_i)^k V_i$ is the same as the nullity of $T - \lambda_i$ on $(T - \lambda_i)^k V$ and equals the number of Jordan blocks of T having eigenvalue λ_i and size greater than k (so for $k = 0$ this gives the number of Jordan blocks).

30. Let λ be an eigenvalue of the linear transformation T on the finite dimensional vector space V over the field F . Let $r_k = \dim_F(T - \lambda)^k V$ be the rank of the linear transformation $(T - \lambda)^k$ on V . For any $k \geq 1$, prove that $r_{k-1} - 2r_k + r_{k+1}$ is the number of Jordan blocks of T corresponding to λ of size k [use Exercise 12 in Section 1]. (This gives an efficient method for determining the Jordan canonical form for T by computing the ranks of the matrices $(A - \lambda I)^k$ for a matrix A representing T , cf. Exercise 31(a) in Section 11.2.)
31. Let N be an $n \times n$ matrix with coefficients in the field F . The matrix N is said to be *nilpotent* if some power of N is the zero matrix, i.e., $N^k = 0$ for some k . Prove that any nilpotent matrix is similar to a block diagonal matrix whose blocks are matrices with 1's along the first superdiagonal and 0's elsewhere.
32. Prove that if N is an $n \times n$ nilpotent matrix then in fact $N^n = 0$.
33. Let A be a strictly upper triangular $n \times n$ matrix (all entries on and below the main diagonal are zero). Prove that A is nilpotent.
34. Prove that the trace of a nilpotent $n \times n$ matrix is 0 (recall the trace of a matrix is the sum of the diagonal elements).
35. For $0 \leq i \leq n$, let d_i be the g.c.d. of the determinants of all the $i \times i$ minors of $xI - A$, for A as in Theorem 21 (take the 0×0 minor to be 1). Prove that the i^{th} element along the diagonal of the Smith Normal Form for A is d_i/d_{i-1} . This gives the invariant factors for A . [Show these g.c.d.s do not change under elementary row and column operations.]
36. Let $V = \mathbb{C}^n$ be the usual n -dimensional vector space of n -tuples $(\alpha_1, \alpha_2, \dots, \alpha_n)$ of complex numbers. Let T be the linear transformation defined by setting $T(\alpha_1, \alpha_2, \dots, \alpha_n)$ equal to $(0, \alpha_1, \alpha_2, \dots, \alpha_{n-1})$. Determine the Jordan canonical form for T .
37. Let J be a Jordan block of size n with eigenvalue λ over \mathbb{C} .
- Prove that the Jordan canonical form for the matrix J^2 is the Jordan block of size n with eigenvalue λ^2 if $\lambda \neq 0$.
 - If $\lambda = 0$ prove that the Jordan canonical form for J^2 has two blocks (with eigenvalues 0) of size $\frac{n}{2}, \frac{n}{2}$ if n is even and of size $\frac{n-1}{2}, \frac{n+1}{2}$ if n is odd.
38. Determine necessary and sufficient conditions for a matrix $A \in M_n(\mathbb{C})$ to have a square root, i.e., for there to exist another matrix $B \in M_n(\mathbb{C})$ such that $A = B^2$. [Suppose B is in Jordan canonical form and consider the Jordan canonical form for B^2 using the previous exercise.]
39. Let J be a Jordan block of size n with eigenvalue λ over a field F of characteristic 2. Determine the Jordan canonical form for the matrix J^2 . Determine necessary and sufficient conditions for a matrix $A \in M_n(F)$ to have a square root, i.e., for there to exist another matrix $B \in M_n(F)$ such that $A = B^2$.

The remaining exercises explore functions (power series) of a matrix and introduce some applications of the Jordan canonical form to the theory of differential equations.

Throughout these exercises the matrices are assumed to be $n \times n$ matrices with entries from the field K , where K is either the real or complex numbers. Let

$$G(x) = \sum_{k=0}^{\infty} \alpha_k x^k$$

be a power series with coefficients from K . Let $G_N(x) = \sum_{k=0}^N \alpha_k x^k$ be the N^{th} partial sum of $G(x)$ and for each $A \in M_n(K)$ let $G_N(A)$ be the element of $M_n(K)$ obtained (as usual) by substituting A in this polynomial. For each fixed i, j we obtain a sequence of real or complex

numbers c_{ij}^N , $N = 0, 1, 2, \dots$ by taking c_{ij}^N to be the i, j entry of the matrix $G_N(A)$. The series

$$G(A) = \sum_{k=0}^{\infty} \alpha_k A^k$$

is said to *converge* to the matrix C in $M_n(K)$ if for each $i, j \in \{1, 2, \dots, n\}$ the sequence c_{ij}^N , $N = 0, 1, 2, \dots$ converges to the i, j entry of C (in which case we write $G(A) = C$). Say $G(A)$ *converges* if there is some $C \in M_n(K)$ such that $G(A) = C$. If A is a 1×1 matrix, this is the usual notion of convergence of a series in K .

For $A = (a_{ij}) \in M_n(K)$ define

$$\|A\| = \sum_{i,j=1}^n |a_{ij}|$$

i.e., $\|A\|$ is the sum of the absolute values of all the entries of A .

- 40.** Prove that for all $A, B \in M_n(K)$ and all $\alpha \in K$

- (a) $\|A + B\| \leq \|A\| + \|B\|$
- (b) $\|AB\| \leq \|A\| \cdot \|B\|$
- (c) $\|\alpha A\| = |\alpha| \cdot \|A\|$.

- 41.** Let R be the radius of convergence of the real or complex power series $G(x)$ (where $R = \infty$ if $G(x)$ converges for all $x \in K$).

- (a) Prove that if $\|A\| < R$ then $G(A)$ converges.

- (b) Deduce that for *all* matrices A the following power series converge:

$$\begin{aligned}\sin(A) &= A - \frac{A^3}{3!} + \frac{A^5}{5!} + \cdots + (-1)^k \frac{A^{2k+1}}{(2k+1)!} + \cdots \\ \cos(A) &= I - \frac{A^2}{2!} + \frac{A^4}{4!} + \cdots + (-1)^k \frac{A^{2k}}{(2k)!} + \cdots \\ \exp(A) &= I + A + \frac{A^2}{2!} + \frac{A^3}{3!} + \cdots + \frac{A^k}{k!} + \cdots\end{aligned}$$

where I is the $n \times n$ identity matrix.

In view of applications to the theory of differential equations we introduce a variable t at this point, so that for $A \in M_n(K)$ the matrix At is obtained from A by multiplying each entry by t (which is the same as multiplying A by the “scalar” matrix tI). We obtain a function from a subset of K into $M_n(K)$ defined by $t \mapsto G(At)$ at all points t where the series $G(At)$ converges. In particular, $\sin(At)$, $\cos(At)$ and $\exp(At)$ converge for all $t \in K$.

- 42.** Let P be a nonsingular $n \times n$ matrix.

- (a) Prove that $PG(At)P^{-1} = G(PAtP^{-1}) = G(PAP^{-1}t)$. (This implies that, up to a change of basis, it suffices to compute $G(At)$ for matrices A in canonical form). [Take limits of partial sums to get the first equality. The second equality is immediate because the matrix tI commutes with every matrix.]
- (b) Prove that if A is the direct sum of matrices A_1, A_2, \dots, A_m , then $G(At)$ is the direct sum of the matrices $G(A_1t), G(A_2t), \dots, G(A_mt)$.
- (c) Show that if Z is the diagonal matrix with entries z_1, z_2, \dots, z_n then $G(Zt)$ is the diagonal matrix with entries $G(z_1t), G(z_2t), \dots, G(z_nt)$.

The matrix $\exp(A)$ defined in Exercise 41(b) is called the *exponential* of A and is often denoted by e^A . The next three exercises lead to a formula for the matrix $\exp(Jt)$, where J is an elementary Jordan matrix.

43. Prove that if A and B are commuting matrices then $\exp(A + B) = \exp(A)\exp(B)$. [Treat A and B as commuting indeterminates and deduce this by comparing the power series on the left hand side with the product of the two power series on the right hand side.]

44. Use the preceding exercise to show that if M is any matrix and λ is any element of K then

$$\exp(\lambda It + M) = e^{\lambda t} \exp(M).$$

45. Let N be the $r \times r$ matrix with 1's on the first superdiagonal and zeros elsewhere. Compute the exponential of the following nilpotent $r \times r$ matrix:

$$\text{if } Nt = \begin{pmatrix} 0 & t & & & \\ & 0 & t & & \\ & & \ddots & & \\ & & & t & \\ & & & & 0 \end{pmatrix} \quad \text{then } \exp(Nt) = \begin{pmatrix} 1 & t & \frac{t^2}{2!} & \cdots & \cdots & \cdots & \frac{t^{r-1}}{(r-1)!} \\ 1 & t & \frac{t^2}{2!} & & & & \vdots \\ \ddots & \ddots & \ddots & \ddots & & & \vdots \\ \ddots & & t & \frac{t^2}{2!} & & & \\ 1 & t & t & & & & \\ & & & & & & 1 \end{pmatrix}.$$

Deduce that if J is the $r \times r$ elementary Jordan matrix with eigenvalue λ then

$$\exp(Jt) = \begin{pmatrix} e^{\lambda t} & te^{\lambda t} & \frac{t^2}{2!}e^{\lambda t} & \cdots & \cdots & \cdots & \frac{t^{r-1}}{(r-1)!}e^{\lambda t} \\ e^{\lambda t} & te^{\lambda t} & \frac{t^2}{2!}e^{\lambda t} & & & & \vdots \\ \ddots & \ddots & \ddots & \ddots & & & \vdots \\ \ddots & & te^{\lambda t} & \frac{t^2}{2!}e^{\lambda t} & & & \\ e^{\lambda t} & & te^{\lambda t} & e^{\lambda t} & te^{\lambda t} & & \\ & & & e^{\lambda t} & & & \end{pmatrix}.$$

[To do the first part use the observation that since Nt is a nilpotent matrix, $\exp(Nt)$ is a polynomial in Nt , i.e., all but a finite number of the terms in the power series are zero. To compute the exponential of Jt write Jt as $\lambda It + Nt$ and use Exercise 44 with $M = Nt$.]

Let $A \in M_n(K)$ and let P be a change of basis matrix such that $P^{-1}AP$ is in Jordan canonical form. Suppose $P^{-1}AP$ is the sum of elementary Jordan matrices J_1, \dots, J_m . The preceding exercises (with $t = 1$) show that $\exp(A)$ can easily be found by writing $E = \exp(P^{-1}AP)$ as the direct sum of the matrices $\exp(J_1), \dots, \exp(J_m)$ and then changing the basis back again to obtain $\exp(A) = PEP^{-1}$.

46. For the 4×4 matrices D and P given in Example 3 of this section:

$$D = \begin{pmatrix} 1 & 2 & -4 & 4 \\ 2 & -1 & 4 & -8 \\ 1 & 0 & 1 & -2 \\ 0 & 1 & -2 & 3 \end{pmatrix} \quad P = \begin{pmatrix} 0 & 1 & 2 & 0 \\ 2 & 0 & -2 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

show that

$$E = \begin{pmatrix} e & e & 0 & 0 \\ 0 & e & 0 & 0 \\ 0 & 0 & e & e \\ 0 & 0 & 0 & e \end{pmatrix} \quad \text{and} \quad \exp(D) = \begin{pmatrix} e & 2e & -4e & 4e \\ 2e & -e & 4e & -8e \\ e & 0 & e & -2e \\ 0 & e & -2e & 3e \end{pmatrix}.$$

- 47.** Compute the exponential of each of the following matrices:
- the matrix A in Example 2 of this section
 - the matrix in Exercise 4 (where you computed the Jordan canonical form and a change of basis matrix)
 - the matrix in Exercise 16.
- 48.** Show that $\exp(0) = I$ (here 0 is the zero matrix and I is the identity matrix). Deduce that $\exp(A)$ is nonsingular with inverse $\exp(-A)$ for all matrices $A \in M_n(K)$.
- 49.** Prove that $\det(\exp(A)) = e^{\text{tr}(A)}$, where $\text{tr}(A)$ is the trace of A (the sum of the diagonal entries of A).
- 50.** Fix any $A \in M_n(K)$. Prove that the map

$$K \rightarrow GL_n(K) \quad \text{defined by} \quad t \mapsto \exp(At)$$

is a group homomorphism (here K is the additive group of the field). (Note how this generalizes the familiar exponential map from K to K^\times , which is the $n = 1$ case. The subgroup $\{\exp(At) \mid t \in K\}$ is called a *1-parameter subgroup* of $GL_n(K)$. These subgroups and the exponential map play an important role in the theory of *Lie groups* — $GL_n(K)$ being a particular example of a Lie group.).

Let $G(x)$ be a power series having an infinite radius of convergence and fix a matrix $A \in M_n(K)$. The entries of the matrix $G(At)$ are K -valued functions of the variable t that are defined for all t . Let $c_{ij}(t)$ be the function of t in the i, j entry of $G(At)$. The *derivative* of $G(At)$ with respect to t , denoted by $\frac{d}{dt}G(At)$, is the matrix whose i, j entry is $\frac{d}{dt}c_{ij}(t)$ obtained by differentiating each of the entries of $G(At)$. In other words, if we identify $M_n(K)$ with K^{n^2} by considering each $n \times n$ matrix as an n^2 -tuple, then $t \mapsto G(At)$ is a map from K to K^{n^2} (i.e., is a vector valued function of t) whose derivative is just the usual (componentwise) derivative of this vector valued function.

- 51.** Establish the following properties of derivatives:

- (a) If $G(x) = \sum_{k=0}^{\infty} \alpha_k x^k$ then $\frac{d}{dt}G(At) = A \sum_{k=1}^{\infty} k \alpha_k (At)^{k-1}$.
- (b) If v is an $n \times 1$ matrix with (constant) entries from K then

$$\frac{d}{dt}(G(At)v) = \left(\frac{d}{dt}G(At) \right) v.$$

- 52.** Deduce from part (a) of the preceding exercise that

$$\frac{d}{dt} \exp(At) = A \exp(At).$$

Now let $y_1(t), \dots, y_n(t)$ be differentiable functions of the real variable t that are related by the following linear system of first order differential equations with constant coefficients $a_{ij} \in K$:

$$\begin{aligned} y'_1 &= a_{11}y_1 + a_{12}y_2 + \dots + a_{1n}y_n \\ y'_2 &= a_{21}y_1 + a_{22}y_2 + \dots + a_{2n}y_n \\ &\vdots \\ y'_n &= a_{n1}y_1 + a_{n2}y_2 + \dots + a_{nn}y_n \end{aligned} \tag{*}$$

(here the primes denote derivatives with respect to t). Let A be the matrix whose i, j entry is a_{ij} , so that $(*)$ may be written as

$$\begin{pmatrix} y'_1 \\ y'_2 \\ \vdots \\ y'_n \end{pmatrix} = A \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

or, more succinctly, as $y' = Ay$, where y is the column vector of functions $y_1(t), \dots, y_n(t)$.

An $n \times n$ matrix whose entries are functions of t and whose columns are independent solutions to the system $(*)$ is called a *fundamental matrix* of $(*)$. By the theory of differential equations, the set of vectors y that are solutions to the system $(*)$ form an n -dimensional vector space over K and so the columns of a fundamental matrix are a *basis for the vector space of all solutions to $(*)$* .

- 53.** Prove that $\exp(At)$ is a fundamental matrix of $(*)$. Show also that if C is the $n \times 1$ constant vector whose entries are $y_1(0), \dots, y_n(0)$ then $y(t) = \exp(At)C$ is the particular solution to the system $(*)$ satisfying the initial condition $y(0) = C$. (Note how this generalizes the 1-dimensional result that the single differential equation $y' = ay$ has e^{at} as a basis for the 1-dimensional space of solutions and the unique solution to this differential equation satisfying the initial condition $y(0) = c$ is $y = ce^{at}$.) [Use the preceding exercises.]
- 54.** Prove that if M is a fundamental matrix of $(*)$ and if Q is a nonsingular matrix in $M_n(K)$, then MQ is also a fundamental matrix of $(*)$. [The columns of MQ are linear combinations of the columns of M .]

Now apply the preceding two exercises to solve some specific systems of differential equations as follows: given the matrix A in a system $(*)$, calculate a change of basis matrix P such that $B = P^{-1}AP$ is in Jordan canonical form. Then $\exp(At) = P \exp(Bt)P^{-1}$ is a fundamental matrix for $(*)$. By the preceding exercise, $P \exp(Bt)$ is also a fundamental matrix for $(*)$ and $\exp(Bt)$ can be calculated by the method described in the discussion following Exercise 45 (in particular, one does not have to find the inverse of the matrix P to obtain a fundamental matrix for $(*)$). Thus, for example, if $A = D$ and P are the matrices given in Exercise 46, then we saw that the Jordan canonical form for A is the matrix $B = P^{-1}AP$ consisting of two 2×2 Jordan blocks with eigenvalues 1. A fundamental matrix for the system $y' = Ay$ is therefore

$$P \exp(Bt) = \begin{pmatrix} 0 & 1 & 2 & 0 \\ 2 & 0 & -2 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} e^t & te^t & 0 & 0 \\ 0 & e^t & 0 & 0 \\ 0 & 0 & e^t & te^t \\ 0 & 0 & 0 & e^t \end{pmatrix} = \begin{pmatrix} 0 & e^t & 2e^t & 2te^t \\ 2e^t & 2te^t & -2e^t & e^t(1-2t) \\ e^t & te^t & 0 & 0 \\ 0 & 0 & e^t & te^t \end{pmatrix}.$$

Writing this out more explicitly, this shows that the general solution to the system of differential equations

$$\begin{aligned} y'_1 &= y_1 + 2y_2 - 4y_3 + 4y_4 \\ y'_2 &= 2y_1 - y_2 + 4y_3 - 8y_4 \\ y'_3 &= y_1 + y_3 - 2y_4 \\ y'_4 &= y_2 - 2y_3 + 3y_4 \end{aligned}$$

is given by

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix} = \alpha_1 \begin{pmatrix} 0 \\ 2e^t \\ e^t \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} e^t \\ 2te^t \\ te^t \\ 0 \end{pmatrix} + \alpha_3 \begin{pmatrix} 2e^t \\ -2e^t \\ 0 \\ e^t \end{pmatrix} + \alpha_4 \begin{pmatrix} 2te^t \\ e^t(1-2t) \\ 0 \\ te^t \end{pmatrix}$$

where $\alpha_1, \dots, \alpha_4$ are arbitrary elements of the field K (this describes the 4-dimensional vector space of solutions).

55. In each of Parts (a) to (c) find a fundamental matrix for the system (*), where the coefficient matrix A of (*) is specified.
- A is the matrix in Part (a) of Exercise 47.
 - A is the matrix in Part (b) of Exercise 47.
 - A is the matrix in Part (c) of Exercise 47.
56. Consider the system (*) whose coefficient matrix A is the matrix D listed in Exercise 46 and whose fundamental matrix was computed just before the preceding exercise. Find the particular solution to (*) that satisfies the initial condition $y_i(0) = 1$ for $i = 1, 2, 3, 4$.

Next we explore a special case of (*). Given the linear n^{th} order differential equation with constant coefficients

$$y^{(n)} + a_{n-1}y^{(n-1)} + \cdots + a_1y' + a_0y = 0 \quad (**)$$

(where $y^{(k)}$ is the k^{th} derivative of y and $y^{(0)} = y$) one can form a *system* of linear *first order* differential equations by letting $y_i = y^{(i-1)}$ for $1 \leq i \leq n$ (the coefficient matrix of this system is described in the next exercise). A basis for the n -dimensional vector space of solutions to the n^{th} order equation (**) may then be obtained from a fundamental matrix for the linear system. Specifically, in each of the $n \times 1$ columns of functions in a fundamental matrix for the system, the $1, 1$ entry is a solution to (**) and so the n functions in the first row of the fundamental matrix for the system form a basis for the solutions to (***).

57. Prove that the matrix, A , of coefficients of the system of n first order equations obtained from (**) is the transpose of the companion matrix of the polynomial $x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0$.
58. Use the above methods to find a basis for the vector space of solutions to the following differential equations
- $y''' - 3y' + 2y = 0$
 - $y'''' + 4y''' + 6y'' + 4y' + y = 0$.

A system of differential equations

$$\begin{aligned} y'_1 &= F_1(y_1, y_2, \dots, y_n) \\ y'_2 &= F_2(y_1, y_2, \dots, y_n) \\ &\vdots \\ y'_n &= F_n(y_1, y_2, \dots, y_n) \end{aligned}$$

where F_1, F_2, \dots, F_n are functions of n variables, is called an *autonomous* system and it will be written more succinctly as $y' = F(y)$, where $F = (F_1, \dots, F_n)$. (The expression autonomous means “independent of time” and it indicates that the variable t — which may be thought of as a time variable — does not appear explicitly on the right hand side.) The system (*) is the special type of autonomous system in which each F_i is a linear function. In many instances it is desirable to analyze the behavior of solutions to an autonomous system of differential equations without explicitly finding these solutions (indeed, it is unlikely that it will be possible to find explicit solutions for a given nonlinear system). This investigation falls under the rubric “qualitative analysis” of autonomous differential equations and the rudiments of this study are often treated in basic calculus courses for 1×1 systems. The first step in a qualitative analysis of an $n \times n$ autonomous system is to find the *steady states*, namely the

constant solutions (these are called steady states since they do not change with t). Note that a constant function $y = c$, where c is the $n \times 1$ constant vector with entries c_1, \dots, c_n , is a solution to $y' = F(y)$ if and only if

$$c'_i = 0 = F_i(c_1, \dots, c_n) \quad \text{for } i = 1, 2, \dots, n,$$

so the steady states are found by computing the zeros of F (in the case of a nonlinear system this may require numerical methods). Next, given the initial value of some solution, one wishes to analyze the behavior of this solution as $t \rightarrow \infty$. This is called the *asymptotic behavior* of the solution. Again, it may not be possible to find the solution explicitly, although by the general theory of differential equations a solution to the initial value problem is unique provided the functions F_i are differentiable. A steady state $y = c$ is called *globally asymptotically stable* if every solution tends to c as $t \rightarrow \infty$, i.e., for any solution $y(t)$ we have $\lim_{t \rightarrow \infty} y_i(t) = c_i$ for all $i = 1, 2, \dots, n$.

In the case of the linear autonomous system (*) the solutions form a vector space, so the only constant solution is the zero solution. The next exercise gives a *sufficient* condition for zero to be globally asymptotically stable and it gives one example of how the behavior of a linear system may be analyzed in terms of the eigenvalues of its coefficient matrix. Nonlinear systems can be approximated by linear systems in some neighborhood of a steady state by considering $y' = Ty$, where $T = \left(\frac{\partial F_i}{\partial y_j} \right)$ is the $n \times n$ Jacobian matrix of F evaluated at the steady state point. In this way the analysis of linear systems plays an important role in the local analysis of general autonomous systems.

- 59.** Prove that the solution of (*) given by $y_i(t) = 0$ for all $i \in \{1, \dots, n\}$ (i.e., the zero solution) is globally asymptotically stable if all the eigenvalues of A have negative real parts. [For those unfamiliar with the behavior of the complex exponential function, assume all eigenvalues are real (hence are negative real numbers). Use the explicit nature of the solutions to show that they all tend to zero as $t \rightarrow \infty$.]

Part IV

FIELD THEORY AND GALOIS THEORY

The previous sections have developed the theory of some of the basic algebraic structures of groups, rings and fields. The next two chapters consider properties of fields, particularly fields which arise from trying to solve equations (such as the simple equation $x^2 + 1 = 0$), and fields which naturally arise in trying to perform “arithmetic” (adding, subtracting, multiplying and dividing). The elegant and beautiful Galois Theory relates the structure of *fields* to certain related *groups* and is one of the basic algebraic tools. Applications include solutions of classical compass and straightedge construction questions, finite fields and Abel’s famous theorem on the insolvability (by radicals) of the general quintic polynomial.

Field Theory

13.1 BASIC THEORY OF FIELD EXTENSIONS

Recall that a field F is a commutative ring with identity in which every nonzero element has an inverse. Equivalently, the set $F^\times = F - \{0\}$ of nonzero elements of F is an abelian group under multiplication.

One of the first invariants associated with any field F is its *characteristic*, defined as follows: If 1_F denotes the identity of F , then F contains the elements $1_F, 1_F + 1_F, 1_F + 1_F + 1_F, \dots$ of the additive subgroup of F generated by 1_F , which may not all be distinct. For n a positive integer, let $n \cdot 1_F = 1_F + \dots + 1_F$ (n times). Then two possibilities arise: either all the elements $n \cdot 1_F$ are distinct, or else $n \cdot 1_F = 0$ for some positive integer n .

Definition. The *characteristic* of a field F , denoted $\text{ch}(F)$, is defined to be the smallest positive integer p such that $p \cdot 1_F = 0$ if such a p exists and is defined to be 0 otherwise.

It is easy to see that

$$\begin{aligned} n \cdot 1_F + m \cdot 1_F &= (m + n) \cdot 1_F && \text{and that} \\ (n \cdot 1_F)(m \cdot 1_F) &= mn \cdot 1_F \end{aligned} \tag{13.1}$$

for positive integers m and n . It follows that the characteristic of a field is either 0 or a prime p (hence the choice of p in the definition above), since if $n = ab$ is composite with $n \cdot 1_F = 0$, then $ab \cdot 1_F = (a \cdot 1_F)(b \cdot 1_F) = 0$ and since F is a field, one of $a \cdot 1_F$ or $b \cdot 1_F$ is 0, so the smallest such integer is necessarily a prime. It also follows that if $n \cdot 1_F = 0$, then n is divisible by p .

Proposition 1. The characteristic of a field F , $\text{ch}(F)$, is either 0 or a prime p . If $\text{ch}(F) = p$ then for any $\alpha \in F$,

$$p \cdot \alpha = \underbrace{\alpha + \alpha + \cdots + \alpha}_{p \text{ times}} = 0.$$

Proof: Only the second statement has not been proved, and this follows immediately from the evident equality $p \cdot \alpha = p \cdot (1_F \alpha) = (p \cdot 1_F)(\alpha)$ in F .

Remark: This notion of a characteristic makes sense also for any integral domain and its characteristic will be the same as for its field of fractions.

Examples

- (1) The fields \mathbb{Q} and \mathbb{R} both have characteristic 0: $\text{ch}(\mathbb{Q}) = \text{ch}(\mathbb{R}) = 0$. The integral domain \mathbb{Z} also has characteristic 0.
- (2) The (finite) field $\mathbb{F}_p = \mathbb{Z}/p\mathbb{Z}$ has characteristic p for any prime p .
- (3) The integral domain $\mathbb{F}_p[x]$ of polynomials in the variable x with coefficients in the field \mathbb{F}_p has characteristic p , as does its field of fractions $\mathbb{F}_p(x)$ (the field of rational functions in x with coefficients in \mathbb{F}_p).

If we define $(-n) \cdot 1_F = -(n \cdot 1_F)$ for positive n and $0 \cdot 1_F = 0$, then we have a natural ring homomorphism (by equation (1))

$$\begin{aligned}\varphi : \mathbb{Z} &\longrightarrow F \\ n &\longmapsto n \cdot 1_F\end{aligned}$$

and we can interpret the characteristic of F by noting that $\ker(\varphi) = \text{ch}(F)\mathbb{Z}$. Taking the quotient by the kernel gives us an *injection* of either \mathbb{Z} or $\mathbb{Z}/p\mathbb{Z}$ into F (depending on whether $\text{ch}(F) = 0$ or $\text{ch}(F) = p$). Since F is a field, we see that F contains a subfield isomorphic either to \mathbb{Q} (the field of fractions of \mathbb{Z}) or to $\mathbb{F}_p = \mathbb{Z}/p\mathbb{Z}$ (the field of fractions of $\mathbb{Z}/p\mathbb{Z}$) depending on the characteristic of F , and in either case is the smallest subfield of F containing 1_F (the field *generated* by 1_F in F).

Definition. The *prime subfield* of a field F is the subfield of F generated by the multiplicative identity 1_F of F . It is (isomorphic to) either \mathbb{Q} (if $\text{ch}(F) = 0$) or \mathbb{F}_p (if $\text{ch}(F) = p$).

Remark: We shall usually denote the identity 1_F of a field F simply by 1. Then in a field of characteristic p , one has $p \cdot 1 = 0$, frequently written simply $p = 0$ (for example, 2 = 0 in a field of characteristic 2). It should be kept in mind, however, that this is a shorthand statement — the element “ p ” is really $p \cdot 1_F$ and is not a distinct element in F . This notation is useful in light of the second statement in Proposition 1.

Examples

- (1) The prime subfield of both \mathbb{Q} and \mathbb{R} is \mathbb{Q} .
- (2) The prime subfield of the field $\mathbb{F}_p(x)$ is isomorphic to \mathbb{F}_p , given by the constant polynomials.

Definition. If K is a field containing the subfield F , then K is said to be an *extension field* (or simply an *extension*) of F , denoted K/F or by the diagram

$$\begin{array}{c} K \\ \downarrow \\ F \end{array}$$

In particular, every field F is an extension of its prime subfield. The field F is sometimes called the *base field* of the extension.

The notation K/F for a field extension is a shorthand for “ K over F ” and is not the quotient of K by F .

If K/F is any extension of fields, then the multiplication defined in K makes K into a *vector space* over F . In particular every field F can be considered as a vector space over its prime field.

Definition. The *degree* (or *relative degree* or *index*) of a field extension K/F , denoted $[K : F]$, is the dimension of K as a vector space over F (i.e., $[K : F] = \dim_F K$). The extension is said to be *finite* if $[K : F]$ is finite and is said to be *infinite* otherwise.

An important class of field extensions are those obtained by trying to solve equations over a given field F . For example, if $F = \mathbb{R}$ is the field of real numbers, then the simple equation $x^2 + 1 = 0$ does not have a solution in F . The question arises whether there is some larger field containing \mathbb{R} in which this equation does have a solution, and it was this question that led Gauss to introduce the *complex numbers* $\mathbb{C} = \mathbb{R} + \mathbb{R}i$, where i is defined so that $i^2 + 1 = 0$. One then defines addition and multiplication in \mathbb{C} by the usual rules familiar from elementary algebra and checks that in fact \mathbb{C} so defined is a *field*, i.e., it is possible to find an inverse for every nonzero element of \mathbb{C} .

Given any field F and any polynomial $p(x) \in F[x]$ one can ask a similar question: does there exist an extension K of F containing a solution of the equation $p(x) = 0$ (i.e., containing a *root* of $p(x)$)? Note that we may assume here that the polynomial $p(x)$ is irreducible in $F[x]$ since a root of any factor of $p(x)$ is certainly a root of $p(x)$ itself. The answer is yes and follows almost immediately from our work on the polynomial ring $F[x]$. We first recall the following useful result on homomorphisms of fields (Corollary 10 of Chapter 7) which follows from the fact that the only ideals of a field F are 0 and F .

Proposition 2. Let $\varphi : F \rightarrow F'$ be a homomorphism of fields. Then φ is either identically 0 or is injective, so that the image of φ is either 0 or isomorphic to F .

Theorem 3. Let F be a field and let $p(x) \in F[x]$ be an irreducible polynomial. Then there exists a field K containing an isomorphic copy of F in which $p(x)$ has a root. Identifying F with this isomorphic copy shows that there exists an extension of F in which $p(x)$ has a root.

Proof: Consider the quotient

$$K = F[x]/(p(x))$$

of the polynomial ring $F[x]$ by the ideal generated by $p(x)$. Since by assumption $p(x)$ is an irreducible polynomial in the P.I.D. $F[x]$, the ideal $(p(x))$ is a *maximal* ideal. Hence K is actually a *field* (this is Proposition 12 of Chapter 7). The canonical projection π of $F[x]$ to the quotient $F[x]/(p(x))$ restricted to $F \subset F[x]$ gives a homomorphism $\varphi = \pi|_F : F \rightarrow K$ which is not identically 0 since it maps the identity 1 of F to the identity 1 of K . Hence by the proposition above, $\varphi(F) \cong F$ is an isomorphic copy

of F contained in K . We identify F with its isomorphic image in K and view F as a subfield of K . If $\bar{x} = \pi(x)$ denotes the image of x in the quotient K , then

$$\begin{aligned} p(\bar{x}) &= \overline{p(x)} && (\text{since } \pi \text{ is a homomorphism}) \\ &= p(x) \pmod{p(x)} && \text{in } F[x]/(p(x)) \\ &= 0 && \text{in } F[x]/(p(x)) \end{aligned}$$

so that K does indeed contain a root of the polynomial $p(x)$. Then K is an extension of F in which the polynomial $p(x)$ has a root.

We shall use this result later to construct extensions of F containing *all* the roots of $p(x)$ (this is the notion of a *splitting field* and one of the central objects of interest in Galois theory).

To understand the field $K = F[x]/(p(x))$ constructed above more fully, it is useful to have a simple representation for the elements of this field. Since F is a subfield of K , we might in particular ask for a basis for K as a vector space over F .

Theorem 4. Let $p(x) \in F[x]$ be an irreducible polynomial of degree n over the field F and let K be the field $F[x]/(p(x))$. Let $\theta = x \pmod{p(x)} \in K$. Then the elements

$$1, \theta, \theta^2, \dots, \theta^{n-1}$$

are a basis for K as a vector space over F , so the degree of the extension is n , i.e., $[K : F] = n$. Hence

$$K = \{a_0 + a_1\theta + a_2\theta^2 + \dots + a_{n-1}\theta^{n-1} \mid a_0, a_1, \dots, a_{n-1} \in F\}$$

consists of all polynomials of degree $< n$ in θ .

Proof: Let $a(x) \in F[x]$ be any polynomial with coefficients in F . Since $F[x]$ is a Euclidean Domain (this is Theorem 3 of Chapter 9), we may divide $a(x)$ by $p(x)$:

$$a(x) = q(x)p(x) + r(x) \quad q(x), r(x) \in F[x] \text{ with } \deg r(x) < n.$$

Since $q(x)p(x)$ lies in the ideal $(p(x))$, it follows that $a(x) \equiv r(x) \pmod{p(x)}$, which shows that every residue class in $F[x]/(p(x))$ is represented by a polynomial of degree less than n . Hence the images $1, \theta, \theta^2, \dots, \theta^{n-1}$ of $1, x, x^2, \dots, x^{n-1}$ in the quotient span the quotient as a vector space over F . It remains to see that these elements are linearly independent, so form a *basis* for the quotient over F .

If the elements $1, \theta, \theta^2, \dots, \theta^{n-1}$ were not linearly independent in K , then there would be a linear combination

$$b_0 + b_1\theta + b_2\theta^2 + \dots + b_{n-1}\theta^{n-1} = 0$$

in K , with $b_0, b_1, \dots, b_{n-1} \in F$, not all 0. This is equivalent to

$$b_0 + b_1x + b_2x^2 + \dots + b_{n-1}x^{n-1} \equiv 0 \pmod{p(x)}$$

i.e.,

$$p(x) \text{ divides } b_0 + b_1x + b_2x^2 + \dots + b_{n-1}x^{n-1}$$

in $F[x]$. But this is impossible, since $p(x)$ is of degree n and the degree of the nonzero polynomial on the right is $< n$. This proves that $1, \theta, \theta^2, \dots, \theta^{n-1}$ are a basis for K over F , so that $[K : F] = n$ by definition. The last statement of the theorem is clear.

This theorem provides an easy description of the elements of the field $F[x]/(p(x))$ as polynomials of degree $< n$ in θ where θ is an element (in K) with $p(\theta) = 0$. It remains only to see how to add and multiply elements written in this form. The addition in the quotient $F[x]/(p(x))$ is just usual addition of polynomials. The multiplication of polynomials $a(x)$ and $b(x)$ in the quotient $F[x]/(p(x))$ is performed by finding the product $a(x)b(x)$ in $F[x]$, then finding the representative of degree $< n$ for the coset $a(x)b(x) + (p(x))$ (as in the proof above) by dividing $a(x)b(x)$ by $p(x)$ and finding the remainder.

This can also be done easily in terms of θ as follows: We may suppose $p(x)$ is monic (since its roots and the ideal it generates do not change by multiplying by a constant), say $p(x) = x^n + p_{n-1}x^{n-1} + \dots + p_1x + p_0$. Then in K , since $p(\theta) = 0$, we have

$$\theta^n = -(p_{n-1}\theta^{n-1} + \dots + p_1\theta + p_0)$$

i.e., θ^n is a linear combination of lower powers of θ . Multiplying both sides by θ and replacing the θ^n on the right hand side by these lower powers again, we see that also θ^{n+1} is a polynomial of degree $< n$ in θ . Similarly, any positive power of θ can be written as a polynomial of degree $< n$ in θ , hence *any* polynomial in θ can be written as a polynomial of degree $< n$ in θ . Multiplication in K is now easily performed: one simply writes the product of two polynomials of degree $< n$ in θ as another polynomial of degree $< n$ in θ .

We summarize this as:

Corollary 5. Let K be as in Theorem 4, and let $a(\theta), b(\theta) \in K$ be two polynomials of degree $< n$ in θ . Then addition in K is defined simply by usual polynomial addition and multiplication in K is defined by

$$a(\theta)b(\theta) = r(\theta)$$

where $r(x)$ is the remainder (of degree $< n$) obtained after dividing the polynomial $a(x)b(x)$ by $p(x)$ in $F[x]$.

By the results proved above, this definition of addition and multiplication on the polynomials of degree $< n$ in θ make K into a *field*, so that one can also *divide* by nonzero elements as well, which is not so immediately obvious from the definitions of the operations.

It is also important in Theorem 4 that the polynomial $p(x)$ be *irreducible* over F . In general the addition and multiplication in Corollary 5 (which can be defined in the same way for any polynomial $p(x)$) do *not* make the polynomials of degree $< n$ in θ into a field if $p(x)$ is not irreducible. In fact, this set is not even an integral domain in general (its structure is given by Proposition 16 of Chapter 9). To describe the *field* containing a root θ of a general polynomial $f(x)$ over F , $f(x)$ is factored into irreducibles in $F[x]$ and the results above are applied to an irreducible factor $p(x)$ of $f(x)$ having θ as a root. We shall consider this more in the following sections.

Examples

- (1) If we apply this construction to the special case $F = \mathbb{R}$ and $p(x) = x^2 + 1$ then we obtain the field

$$\mathbb{R}[x]/(x^2 + 1)$$

which is an extension of degree 2 of \mathbb{R} in which $x^2 + 1$ has a root. The elements of this field are of the form $a + b\theta$ for $a, b \in \mathbb{R}$. Addition is defined by

$$(a + b\theta) + (c + d\theta) = (a + c) + (b + d)\theta. \quad (13.2a)$$

To multiply we use the fact that $\theta^2 + 1 = 0$, i.e., $\theta^2 = -1$ in K . (Alternatively, note that -1 is also the remainder when x^2 is divided by $x^2 + 1$ in $\mathbb{R}[x]$.) Then

$$\begin{aligned} (a + b\theta)(c + d\theta) &= ac + (ad + bc)\theta + bd\theta^2 \\ &= ac + (ad + bc)\theta + bd(-1) \\ &= (ac - bd) + (ad + bc)\theta. \end{aligned} \quad (13.2b)$$

These are, up to changing θ to i , the formulas for adding and multiplying in \mathbb{C} . Put another way, the map

$$\begin{aligned} \varphi : \mathbb{R}[x]/(x^2 + 1) &\longrightarrow \mathbb{C} \\ a + bx &\mapsto a + bi \end{aligned}$$

is a homomorphism. Since it is bijective (as a map of vector spaces over the reals, for example), it is an isomorphism. Notice that instead of taking the existence of \mathbb{C} for granted (along with the fairly tedious verification that it is in fact a field), we could have *defined* \mathbb{C} by this isomorphism. Then the fact that it is a field is a consequence of Theorem 4.

- (2) Take now $F = \mathbb{Q}$ to be the field of rational numbers and again take $p(x) = x^2 + 1$ (still irreducible over \mathbb{Q} , of course). Then the same construction, with the same addition and multiplication formulas as (2a) and (2b) above, except that now a and b are elements of \mathbb{Q} , defines a field extension $\mathbb{Q}(i)$ of \mathbb{Q} of degree 2 containing a root i of $x^2 + 1$.
- (3) Take $F = \mathbb{Q}$ and $p(x) = x^2 - 2$, irreducible over \mathbb{Q} by Eisenstein's Criterion, for example. Then we obtain a field extension of \mathbb{Q} of degree 2 containing a square root θ of 2, denoted $\mathbb{Q}(\theta)$. If we denote θ by $\sqrt{2}$, the elements of this field are of the form

$$a + b\sqrt{2}, \quad a, b \in \mathbb{Q}$$

with addition defined by

$$(a + b\sqrt{2}) + (c + d\sqrt{2}) = (a + c) + (b + d)\sqrt{2}$$

and multiplication defined by

$$(a + b\sqrt{2})(c + d\sqrt{2}) = (ac + 2bd) + (ad + bc)\sqrt{2}.$$

- (4) Let $F = \mathbb{Q}$ and $p(x) = x^3 - 2$, irreducible again by Eisenstein. Denoting a root of $p(x)$ by θ , we obtain the field

$$\mathbb{Q}[x]/(x^3 - 2) \cong \{a + b\theta + c\theta^2 \mid a, b, c \in \mathbb{Q}\}$$

with $\theta^3 = 2$, an extension of degree 3. To find the inverse of, say, $1 + \theta$ in this field, we can proceed as follows: By the Euclidean Algorithm in $\mathbb{Q}[x]$ there are polynomials $a(x)$ and $b(x)$ with

$$a(x)(1 + x) + b(x)(x^3 - 2) = 1$$

(since $p(x) = x^3 - 2$ is irreducible, it is relatively prime to every polynomial of smaller degree). In the quotient field this equation implies that $a(\theta)$ is the inverse of $1 + \theta$. In this case, a simple computation shows that we can take $a(x) = \frac{1}{3}(x^2 - x + 1)$ (and $b(x) = -\frac{1}{3}$), so that

$$(1 + \theta)^{-1} = \frac{\theta^2 - \theta + 1}{3}.$$

- (5) In general, if $\theta \in K$ is a root of the irreducible polynomial

$$p(x) = p_n x^n + p_{n-1} x^{n-1} + \cdots + p_1 x + p_0$$

we can compute $\theta^{-1} \in K$ from

$$\theta(p_n \theta^{n-1} + p_{n-1} \theta^{n-2} + \cdots + p_1) = -p_0$$

namely

$$\theta^{-1} = \frac{-1}{p_0} (p_n \theta^{n-1} + p_{n-1} \theta^{n-2} + \cdots + p_1) \in K$$

(note that $p_0 \neq 0$ since $p(x)$ is irreducible).

Remark: Determining inverses in extensions of this type may be familiar from elementary algebra in the case of \mathbb{C} or Example 3 under the name “rationalizing denominators.” The last two examples indicates a procedure which is much more general than the ad hoc procedures of elementary algebra.

- (6) Take $F = \mathbb{F}_2$, the finite field with two elements, and $p(x) = x^2 + x + 1$, which we have previously checked is irreducible over \mathbb{F}_2 . Here we obtain a degree 2 extension of \mathbb{F}_2

$$\mathbb{F}_2[x]/(x^2 + x + 1) \cong \{a + b\theta \mid a, b \in \mathbb{F}_2\}$$

where $\theta^2 = -\theta - 1 = \theta + 1$. Multiplication in this field $\mathbb{F}_2(\theta)$ (which contains four elements) is defined by

$$\begin{aligned} (a + b\theta)(c + d\theta) &= ac + (ad + bc)\theta + bd\theta^2 \\ &= ac + (ad + bc)\theta + bd(\theta + 1) \\ &= (ac + bd) + (ad + bc + bd)\theta. \end{aligned}$$

- (7) Let $F = k(t)$ be the field of rational functions in the variable t over a field k (for example, $k = \mathbb{Q}$ or $k = \mathbb{F}_p$). Let $p(x) = x^2 - t \in F[x]$. Then $p(x)$ is irreducible (it is Eisenstein at the prime (t) in $k[t]$). If we denote a root by θ , the corresponding degree 2 field extension $F(\theta)$ consists of the elements

$$\{a(t) + b(t)\theta \mid a(t), b(t) \in F\}$$

where the coefficients $a(t)$ and $b(t)$ are rational functions in t with coefficients in k and where $\theta^2 = t$.

Suppose F is a subfield of a field K and $\alpha \in K$ is an element of K . Then the collection of subfields of K containing both F and α is nonempty (K is such a field, for example). Since the intersection of subfields is again a subfield, it follows that there is a unique minimal subfield of K containing both F and α (the intersection of all subfields with this property). Similar remarks apply if α is replaced by a collection α, β, \dots of elements of K .

Definition. Let K be an extension of the field F and let $\alpha, \beta, \dots \in K$ be a collection of elements of K . Then the smallest subfield of K containing both F and the elements α, β, \dots , denoted $F(\alpha, \beta, \dots)$ is called the field *generated by α, β, \dots over F* .

Definition. If the field K is generated by a single element α over F , $K = F(\alpha)$, then K is said to be a *simple* extension of F and the element α is called a *primitive element* for the extension.

We shall later characterize which extensions of a field F are simple. In particular we shall prove that every finite extension of a field of characteristic 0 is a simple extension.

The connection between the simple extension $F(\alpha)$ generated by α over F where α is a root of some irreducible polynomial $p(x)$ and the field constructed in Theorem 3 is provided by the following:

Theorem 6. Let F be a field and let $p(x) \in F[x]$ be an irreducible polynomial. Suppose K is an extension field of F containing a root α of $p(x)$: $p(\alpha) = 0$. Let $F(\alpha)$ denote the subfield of K generated over F by α . Then

$$F(\alpha) \cong F[x]/(p(x)).$$

Remark: This theorem says that *any* field over F in which $p(x)$ contains a root contains a subfield isomorphic to the extension of F constructed in Theorem 3 and that this field is (up to isomorphism) the smallest extension of F containing such a root. The difference between this result and Theorem 3 is that Theorem 6 *assumes* the existence of a root α of $p(x)$ in some field K and the major point of Theorem 3 is *proving* that there exists such an extension field K .

Proof: There is a natural homomorphism

$$\begin{aligned}\varphi : F[x] &\longrightarrow F(\alpha) \subseteq K \\ a(x) &\longmapsto a(\alpha)\end{aligned}$$

obtained by mapping F to F by the identity map and sending x to α and then extending so that the map is a ring homomorphism (i.e., the polynomial $a(x)$ in x maps to the polynomial $a(\alpha)$ in α). Since $p(\alpha) = 0$ by assumption, the element $p(x)$ is in the kernel of φ , so we obtain an induced homomorphism (also denoted φ):

$$\varphi : F[x]/(p(x)) \longrightarrow F(\alpha).$$

But since $p(x)$ is irreducible, the quotient on the left is a *field*, and φ is not the 0 map (it is the identity on F , for example), hence φ is an isomorphism of the field on the left with its image. Since this image is then a subfield of $F(\alpha)$ containing F and containing α , by the definition of $F(\alpha)$ the map must be surjective, proving the theorem.

Combined with Corollary 5, this determines the field $F(\alpha)$ when α is a root of an irreducible polynomial $p(x)$:

Corollary 7. Suppose in Theorem 6 that $p(x)$ is of degree n . Then

$$F(\alpha) = \{a_0 + a_1\alpha + a_2\alpha^2 + \cdots + a_{n-1}\alpha^{n-1} \mid a_0, a_1, \dots, a_{n-1} \in F\} \subseteq K.$$

Describing fields generated by more than one element is more complicated and we shall return to this question in the following section.

Examples

- (1) In Example 3 above, we have determined the field $\mathbb{Q}(\sqrt{2})$ generated over \mathbb{Q} by the element $\sqrt{2} \in \mathbb{R}$, having suggestively denoted the abstract solution θ of the equation $x^2 - 2 = 0$ by the symbol $\sqrt{2}$, which has an independent meaning in the field \mathbb{R} (namely the *positive* square root of 2 in \mathbb{R}).
- (2) The equation $x^2 - 2 = 0$ has another solution in \mathbb{R} , namely $-\sqrt{2}$, the *negative* square root of 2 in \mathbb{R} . The field generated over \mathbb{Q} by this solution consists of the elements $\{a + b(-\sqrt{2}) \mid a, b \in \mathbb{Q}\}$, and is again isomorphic to the field in Example 3 above (hence also isomorphic to the field just considered, the isomorphism given explicitly by $a + b\sqrt{2} \mapsto a - b\sqrt{2}$). As a subset of \mathbb{R} this is the same set of elements as in Example 1.
- (3) Similarly, if we use the symbol $\sqrt[3]{2}$ to denote the (positive) cube root of 2 in \mathbb{R} , then the field generated by $\sqrt[3]{2}$ over \mathbb{Q} in \mathbb{R} consists of the elements

$$\{a + b\sqrt[3]{2} + c(\sqrt[3]{2})^2 \mid a, b, c \in \mathbb{Q}\}$$

and is isomorphic to the field constructed in Example 4 above.

- (4) The equation $x^3 - 2 = 0$ has no further solutions in \mathbb{R} , but there are two additional solutions in \mathbb{C} given by $\sqrt[3]{2}\left(\frac{-1+i\sqrt{3}}{2}\right)$ and $\sqrt[3]{2}\left(\frac{-1-i\sqrt{3}}{2}\right)$ ($\sqrt{3}$ denoting the positive real square root of 3) as can easily be checked. The fields generated by either of these two elements over \mathbb{Q} are subfields of \mathbb{C} (but not of \mathbb{R}) and are both isomorphic to the field constructed in the previous example (and to Example 4 earlier).

As Theorem 6 indicates, the roots of an irreducible polynomial $p(x)$ are *algebraically indistinguishable* in the sense that the fields obtained by adjoining any root of an irreducible polynomial are isomorphic. In the last two examples above, the fields obtained by adjoining one of the three possible (complex) roots of $x^3 - 2 = 0$ to \mathbb{Q} were all algebraically isomorphic. The fields were distinguished not by their algebraic properties, but by whether their elements were *real*, which involves *continuous* operations.

The fact that different roots of the same irreducible polynomial have the same algebraic properties can be extended slightly, as follows:

Let $\varphi : F \xrightarrow{\sim} F'$ be an isomorphism of fields. The map φ induces a ring isomorphism (also denoted φ)

$$\varphi : F[x] \xrightarrow{\sim} F'[x]$$

defined by applying φ to the coefficients of a polynomial in $F[x]$. Let $p(x) \in F[x]$ be an irreducible polynomial and let $p'(x) \in F'[x]$ be the polynomial obtained by applying the map φ to the coefficients of $p(x)$, i.e., the image of $p(x)$ under φ . The isomorphism φ maps the maximal ideal $(p(x))$ to the ideal $(p'(x))$, so this ideal is also

maximal, which shows that $p'(x)$ is also irreducible in $F'[x]$. The following theorem shows that the fields obtained by adjoining a root of $p(x)$ to F and a root of $p'(x)$ to F' have the same algebraic structure (i.e., are isomorphic):

Theorem 8. Let $\varphi : F \xrightarrow{\sim} F'$ be an isomorphism of fields. Let $p(x) \in F[x]$ be an irreducible polynomial and let $p'(x) \in F'[x]$ be the irreducible polynomial obtained by applying the map φ to the coefficients of $p(x)$. Let α be a root of $p(x)$ (in some extension of F) and let β be a root of $p'(x)$ (in some extension of F'). Then there is an isomorphism

$$\begin{aligned}\sigma : F(\alpha) &\xrightarrow{\sim} F'(\beta) \\ \alpha &\mapsto \beta\end{aligned}$$

mapping α to β and extending φ , i.e., such that σ restricted to F is the isomorphism φ .

Proof: As noted above, the isomorphism φ induces a natural isomorphism from $F[x]$ to $F'[x]$ which maps the maximal ideal $(p(x))$ to the maximal ideal $(p'(x))$. Taking the quotients by these ideals, we obtain an isomorphism of fields

$$F[x]/(p(x)) \xrightarrow{\sim} F'[x]/(p'(x)).$$

By Theorem 6 the field on the left is isomorphic to $F(\alpha)$ and by the same theorem the field on the right is isomorphic to $F'(\beta)$. Composing these isomorphisms, we obtain the isomorphism σ . It is clear that the restriction of this isomorphism to F is φ , completing the proof.

This extension theorem will be of considerable use when we consider Galois Theory later. It can be represented pictorially by the diagram

$$\begin{array}{ccc} \sigma : & F(\alpha) & \xrightarrow{\sim} & F'(\beta) \\ & | & & | \\ \varphi : & F & \xrightarrow{\sim} & F' \end{array}$$

EXERCISES

- Show that $p(x) = x^3 + 9x + 6$ is irreducible in $\mathbb{Q}[x]$. Let θ be a root of $p(x)$. Find the inverse of $1 + \theta$ in $\mathbb{Q}(\theta)$.
- Show that $x^3 - 2x - 2$ is irreducible over \mathbb{Q} and let θ be a root. Compute $(1 + \theta)(1 + \theta + \theta^2)$ and $\frac{1 + \theta}{1 + \theta + \theta^2}$ in $\mathbb{Q}(\theta)$.
- Show that $x^3 + x + 1$ is irreducible over \mathbb{F}_2 and let θ be a root. Compute the powers of θ in $\mathbb{F}_2(\theta)$.
- Prove directly that the map $a + b\sqrt{2} \mapsto a - b\sqrt{2}$ is an isomorphism of $\mathbb{Q}(\sqrt{2})$ with itself.
- Suppose α is a rational root of a monic polynomial in $\mathbb{Z}[x]$. Prove that α is an integer.
- Show that if α is a root of $a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$ then $a_n \alpha$ is a root of the monic polynomial $x^n + a_{n-1} x^{n-1} + a_n a_{n-2} x^{n-2} + \cdots + a_n^{n-2} a_1 x + a_n^{n-1} a_0$.
- Prove that $x^3 - nx + 2$ is irreducible for $n \neq -1, 3, 5$.
- Prove that $x^5 - ax - 1 \in \mathbb{Z}[x]$ is irreducible unless $a = 0, 2$ or -1 . The first two correspond to linear factors, the third corresponds to the factorization $(x^2 - x + 1)(x^3 + x^2 - 1)$.

13.2 ALGEBRAIC EXTENSIONS

Let F be a field and let K be an extension of F .

Definition. The element $\alpha \in K$ is said to be *algebraic* over F if α is a root of some nonzero polynomial $f(x) \in F[x]$. If α is not algebraic over F (i.e., is not the root of any nonzero polynomial with coefficients in F) then α is said to be *transcendental* over F . The extension K/F is said to be *algebraic* if every element of K is algebraic over F .

Note that if α is algebraic over a field F then it is algebraic over any extension field L of F (if $f(x)$ having α as a root has coefficients in F then it also has coefficients in L).

Proposition 9. Let α be algebraic over F . Then there is a unique monic irreducible polynomial $m_{\alpha,F}(x) \in F[x]$ which has α as a root. A polynomial $f(x) \in F[x]$ has α as a root if and only if $m_{\alpha,F}(x)$ divides $f(x)$ in $F[x]$.

Proof: Let $g(x) \in F[x]$ be a polynomial of minimal degree having α as a root. Multiplying $g(x)$ by a constant, we may assume $g(x)$ is monic. Suppose $g(x)$ were reducible in $F[x]$, say $g(x) = a(x)b(x)$ with $a(x), b(x) \in F[x]$ both of degree smaller than the degree of $g(x)$. Then $g(\alpha) = a(\alpha)b(\alpha)$ in K , and since K is a field, either $a(\alpha) = 0$ or $b(\alpha) = 0$, contradicting the minimality of the degree of $g(x)$. It follows that $g(x)$ is a monic irreducible polynomial having α as a root. Suppose now that $f(x) \in F[x]$ is any polynomial having α as a root. By the Euclidean Algorithm in $F[x]$ there are polynomials $q(x), r(x) \in F[x]$ such that

$$f(x) = q(x)g(x) + r(x) \quad \text{with } \deg r(x) < \deg g(x).$$

Then $f(\alpha) = q(\alpha)g(\alpha) + r(\alpha)$ in K and since α is a root of both $f(x)$ and $g(x)$, we obtain $r(\alpha) = 0$, which contradicts the minimality of $g(x)$ unless $r(x) = 0$. Hence $g(x)$ divides any polynomial $f(x)$ in $F[x]$ having α as a root and, in particular, would divide any other monic irreducible polynomial in $F[x]$ having α as a root. This proves that $m_{\alpha,F}(x) = g(x)$ is unique and completes the proof of the proposition.

Corollary 10. If L/F is an extension of fields and α is algebraic over both F and L , then $m_{\alpha,L}(x)$ divides $m_{\alpha,F}(x)$ in $L[x]$.

Proof: This is immediate from the second statement in Proposition 9 applied to L , since $m_{\alpha,F}(x)$ is a polynomial in $L[x]$ having α as a root.

Definition. The polynomial $m_{\alpha,F}(x)$ (or just $m_{\alpha}(x)$ if the field F is understood) in Proposition 9 is called the *minimal polynomial* for α over F . The *degree* of $m_{\alpha}(x)$ is called the *degree* of α .

Note that by the proposition, a monic polynomial over F with α as a root is the minimal polynomial for α over F if and only if it is irreducible over F . Exercise 20

gives one method for computing the minimal polynomial for α over F , and the theory of Gröbner bases can be used to compute the minimal polynomial for other elements in $F(\alpha)$ (cf. Proposition 10 and Exercise 48 in Section 15.1).

Proposition 11. Let α be algebraic over the field F and let $F(\alpha)$ be the field generated by α over F . Then

$$F(\alpha) \cong F[x]/(m_\alpha(x))$$

so that in particular

$$[F(\alpha) : F] = \deg m_\alpha(x) = \deg \alpha,$$

i.e., the degree of α over F is the degree of the extension it generates over F .

Proof: This follows immediately from Theorem 6.

Examples

- (1) The minimal polynomial for $\sqrt{2}$ over \mathbb{Q} is $x^2 - 2$ and $\sqrt{2}$ is of degree 2 over \mathbb{Q} : $[\mathbb{Q}(\sqrt{2}) : \mathbb{Q}] = 2$.
- (2) The minimal polynomial for $\sqrt[3]{2}$ over \mathbb{Q} is $x^3 - 2$ and $\sqrt[3]{2}$ is of degree 3 over \mathbb{Q} : $[\mathbb{Q}(\sqrt[3]{2}) : \mathbb{Q}] = 3$.
- (3) Similarly, for any $n > 1$, the polynomial $x^n - 2$ is irreducible over \mathbb{Q} since it is Eisenstein. Denoting a root of this polynomial by $\sqrt[n]{2}$ (where as usual we reserve this symbol to denote the *positive* n^{th} root of 2 if we want to view this root as an element of \mathbb{R} , and where the symbol denotes any one of the algebraically indistinguishable abstract solutions in general), we have $[\mathbb{Q}(\sqrt[n]{2}) : \mathbb{Q}] = n$.
- (4) The minimal polynomial and the degree of an element α depend on the base field. For example, over \mathbb{R} , the element $\sqrt[3]{2}$ is of degree *one*, with minimal polynomial $m_{\sqrt[3]{2}, \mathbb{R}}(x) = x - \sqrt[3]{2}$.
- (5) Consider the polynomial $p(x) = x^3 - 3x - 1$ over \mathbb{Q} , which is irreducible over \mathbb{Q} since it is a cubic which has no rational root (cf. Proposition 11 of Chapter 9). Hence $[\mathbb{Q}(\alpha) : \mathbb{Q}] = 3$ for any root α of $p(x)$. For future reference we note that a quick sketch of the graph of this function over the real numbers shows that the graph crosses the x -axis precisely once in the interval $[0, 2]$, i.e., there is precisely one real number α , $0 < \alpha < 2$ satisfying $\alpha^3 - 3\alpha - 1 = 0$.

Proposition 12. The element α is algebraic over F if and only if the simple extension $F(\alpha)/F$ is finite. More precisely, if α is an element of an extension of degree n over F then α satisfies a polynomial of degree at most n over F and if α satisfies a polynomial of degree n over F then the degree of $F(\alpha)$ over F is at most n .

Proof: If α is algebraic over F , then the degree of the extension $F(\alpha)/F$ is the degree of the minimal polynomial for α over F . Hence the extension is finite, of degree $\leq n$ if α satisfies a polynomial of degree n . Conversely, suppose α is an element of an extension of degree n over F (for example, if $[F(\alpha) : F] = n$). Then the $n + 1$ elements

$$1, \alpha, \alpha^2, \dots, \alpha^n$$

of $F(\alpha)$ are linearly dependent over F , say

$$b_0 + b_1\alpha + b_2\alpha^2 + \cdots + b_n\alpha^n = 0$$

with $b_0, b_1, b_2, \dots, b_n \in F$ not all 0. Hence α is the root of a nonzero polynomial with coefficients in F (of degree $\leq n$), which proves α is algebraic over F and also proves the second statement of the proposition.

Corollary 13. If the extension K/F is finite, then it is algebraic.

Proof: If $\alpha \in K$, then the subfield $F(\alpha)$ is in particular a subspace of the vector space K over F . Hence $[F(\alpha) : F] \leq [K : F]$ and so α is algebraic over F by the proposition.

Remark: We shall prove below a sort of converse to this result (Theorem 17), but note that there are infinite algebraic extensions (we shall have an example later), so the literal converse of this corollary is not true.

Example: (Quadratic Extensions over Fields of Characteristic $\neq 2$)

Let F be a field of characteristic $\neq 2$ (for example, any field of characteristic 0, such as \mathbb{Q}) and let K be an extension of F of degree 2, $[K : F] = 2$. Let α be any element of K not contained in F . By the proposition above, α satisfies an equation of degree at most 2 over F . This equation cannot be of degree 1, since α is not an element of F by assumption. It follows that the minimal polynomial of α is a monic quadratic

$$m_\alpha(x) = x^2 + bx + c \quad b, c \in F.$$

Since $F \subset F(\alpha) \subseteq K$ and $F(\alpha)$ is already a vector space over F of dimension 2, we have $K = F(\alpha)$.

The roots of this quadratic equation can be determined by the quadratic formula, which is valid over any field of characteristic $\neq 2$ (the formula is obtained as in elementary algebra by completing the square):

$$\alpha = \frac{-b \pm \sqrt{b^2 - 4c}}{2}$$

(the reason for requiring the characteristic of F not be 2 is that we must divide by 2). Here $b^2 - 4c$ is not a square in F since α is not an element of F and the symbol $\sqrt{b^2 - 4c}$ denotes a root of the equation $x^2 - (b^2 - 4c) = 0$ in K (see the end of the next paragraph). Note that here there is no natural choice of one of the roots analogous to choosing the *positive* square root of 2 in \mathbb{R} — the roots are algebraically indistinguishable.

Now $F(\alpha) = F(\sqrt{b^2 - 4c})$ as follows: by the formula above, α is an element of the field on the right, hence $F(\alpha) \subseteq F(\sqrt{b^2 - 4c})$. Conversely, $\sqrt{b^2 - 4c} = \mp(b + 2\alpha)$ shows that $\sqrt{b^2 - 4c}$ is an element of $F(\alpha)$, which gives the reverse inclusion $F(\sqrt{b^2 - 4c}) \subseteq F(\alpha)$ (and incidentally shows that the equation $x^2 - (b^2 - 4c) = 0$ does have a solution in K).

It follows that any extension K of F of degree 2 is of the form $F(\sqrt{D})$ where D is an element of F which is not a square in F , and conversely, every such extension is an extension of degree 2 of F . For this reason, extensions of degree 2 of a field F are called *quadratic* extensions of F .

Suppose that F is a subfield of a field K which in turn is a subfield of a field L . Then there are three associated extension degrees — the dimension of K and L as vector spaces over F , and the dimension of L as a vector space over K .

Theorem 14. Let $F \subseteq K \subseteq L$ be fields. Then

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

i.e. extension degrees are multiplicative, where if one side of the equation is infinite, the other side is also infinite. Pictorially,

$$\begin{array}{ccccc} & & [L:F] & & \\ \overbrace{F & \subseteq & K & \subseteq & L} & & & & \\ & [K:F] & & & [L:K] \end{array}$$

Proof: Suppose first that $[L : K] = m$ and $[K : F] = n$ are finite. Let $\alpha_1, \alpha_2, \dots, \alpha_m$ be a basis for L over K and let $\beta_1, \beta_2, \dots, \beta_n$ be a basis for K over F . Then every element of L can be written as a linear combination

$$a_1\alpha_1 + a_2\alpha_2 + \cdots + a_m\alpha_m$$

where a_1, \dots, a_m are elements of K , hence are F -linear combinations of β_1, \dots, β_n :

$$a_i = b_{i1}\beta_1 + b_{i2}\beta_2 + \cdots + b_{in}\beta_n \quad i = 1, 2, \dots, m \quad (13.3)$$

where the b_{ij} are elements of F . Substituting these expressions in for the coefficients a_i above, we see that every element of L can be written as a linear combination

$$\sum_{\substack{i=1,2,\dots,m \\ j=1,2,\dots,n}} b_{ij}\alpha_i\beta_j$$

of the mn elements $\alpha_i\beta_j$ with coefficients in F . Hence these elements span L as a vector space over F .

Suppose now that we had a linear relation in L

$$\sum_{\substack{i=1,2,\dots,m \\ j=1,2,\dots,n}} b_{ij}\alpha_i\beta_j = 0$$

with coefficients b_{ij} in F . Then defining the elements $a_i \in K$ by equation (3) above, this linear relation could be written

$$a_1\alpha_1 + a_2\alpha_2 + \cdots + a_m\alpha_m = 0.$$

Since the α_i are a basis for L over K , it follows that all the coefficients $a_i, i = 1, 2, \dots, m$ must be 0, i.e., that

$$b_{i1}\beta_1 + b_{i2}\beta_2 + \cdots + b_{in}\beta_n = 0 \quad i = 1, 2, \dots, m$$

in K . Since now the $\beta_j, j = 1, 2, \dots, n$ form a basis for K over F , this implies $b_{ij} = 0$ for all i and j . Hence the elements $\alpha_i\beta_j$ are linearly independent over F , so form a basis for L over F and $[L : F] = mn = [L : K][K : F]$, as claimed.

If $[K : F]$ is infinite, then there are infinitely many elements of K , hence of L , which are linearly independent over F , so that $[L : F]$ is also infinite. Similarly, if $[L : K]$ is infinite, there are infinitely many elements of L linearly independent over K , so certainly linearly independent over F , so again $[L : F]$ is infinite. Finally, if $[L : K]$ and $[K : F]$ are both finite, then the proof above shows $[L : F]$ is finite, so that $[L : F]$ infinite implies at least one of $[L : K]$ and $[K : F]$ is infinite, completing the proof.

Remark: Note the similarity of this result with the result on group orders proved in Part I. As with diagrams involving groups we shall frequently indicate the relative degrees of extensions in field diagrams.

The multiplicativity of extension degrees is extremely useful in computations. A particular application is the following:

Corollary 15. Suppose L/F is a finite extension and let K be any subfield of L containing F , $F \subseteq K \subseteq L$. Then $[K : F]$ divides $[L : F]$.

Proof: This is immediate.

Examples

- (1) The element $\sqrt{2}$ is not contained in the field $\mathbb{Q}(\alpha)$ where α is the real root of $x^3 - 3x - 1$ between 0 and 2, since we have already determined that $[\mathbb{Q}(\sqrt{2}) : \mathbb{Q}] = 2$ and $[\mathbb{Q}(\alpha) : \mathbb{Q}] = 3$ and 2 does not divide 3. Note that it is not so easy to prove directly that $\sqrt{2}$ cannot be written as a rational linear combination of 1, α , α^2 .
- (2) Let as usual $\sqrt[6]{2}$ denote the positive real 6th root of 2. Then $[\mathbb{Q}(\sqrt[6]{2}) : \mathbb{Q}] = 6$. Since $(\sqrt[6]{2})^3 = \sqrt{2}$ we have $\mathbb{Q}(\sqrt{2}) \subset \mathbb{Q}(\sqrt[6]{2})$ and by the multiplicativity of extension degrees, $[\mathbb{Q}(\sqrt[6]{2}) : \mathbb{Q}(\sqrt{2})] = 3$. This gives us the field diagram

$$\begin{array}{ccccccc} & & & 6 & & & \\ & & & \overbrace{\quad\quad\quad} & & & \\ \mathbb{Q} & \subset & \mathbb{Q}(\sqrt{2}) & \subset & \mathbb{Q}(\sqrt[6]{2}) & & \\ & \underbrace{\quad\quad\quad}_{2} & \underbrace{\quad\quad\quad}_{3} & & & & \end{array}$$

In particular, this shows that the minimal polynomial for $\sqrt[6]{2}$ over $\mathbb{Q}(\sqrt{2})$ is of degree 3. It is therefore the polynomial $x^3 - \sqrt{2}$. Note that showing directly that this polynomial is irreducible over $\mathbb{Q}(\sqrt{2})$ is not completely trivial.

By Theorem 14 a finite extension of a finite extension is finite. The next results use this to show that an extension generated by a finite number of algebraic elements is finite (extending Proposition 12).

Definition. An extension K/F is *finitely generated* if there are elements $\alpha_1, \alpha_2, \dots, \alpha_k$ in K such that $K = F(\alpha_1, \alpha_2, \dots, \alpha_k)$.

Recall that the field generated over F by a collection of elements in a field K is the smallest subfield of K containing these elements and F . The next lemma will show that for finitely generated extensions this field can be obtained recursively by a series of simple extensions.

Lemma 16. $F(\alpha, \beta) = (F(\alpha))(\beta)$, i.e., the field generated over F by α and β is the field generated by β over the field $F(\alpha)$ generated by α .

Proof: This follows by the minimality of the fields in question. The field $F(\alpha, \beta)$ contains F and α , hence contains the field $F(\alpha)$, and since it also contains β , we have the inclusion $(F(\alpha))(\beta) \subseteq F(\alpha, \beta)$ by the minimality of the field $(F(\alpha))(\beta)$. Since the field $(F(\alpha))(\beta)$ contains F , α and β , by the minimality of $F(\alpha, \beta)$ we have the reverse inclusion $F(\alpha, \beta) \subseteq (F(\alpha))(\beta)$, which proves the lemma.

By the lemma we have

$$K = F(\alpha_1, \alpha_2, \dots, \alpha_k) = (F(\alpha_1, \alpha_2, \dots, \alpha_{k-1}))(\alpha_k)$$

and so by iterating, we see that K is obtained by taking the field F_1 generated over F by α_1 , then the field F_2 generated over F_1 (this is important) by α_2 , and so on, with $F_k = K$. This gives a sequence of fields:

$$F = F_0 \subseteq F_1 \subseteq F_2 \subseteq \cdots \subseteq F_k = K$$

where

$$F_{i+1} = F_i(\alpha_{i+1}) \quad i = 0, 1, \dots, k-1.$$

Suppose now that the elements $\alpha_1, \alpha_2, \dots, \alpha_k$ are algebraic over F of degrees n_1, n_2, \dots, n_k (so a priori are algebraic over any extension of F). Then the extensions in this sequence are simple extensions of the type considered in Proposition 11. The relative extension degree $[F_{i+1} : F_i]$ is equal to the degree of the minimal polynomial of α_{i+1} over F_i , which is at most n_{i+1} (and equals n_{i+1} if and only if the minimal polynomial of α_{i+1} over F remains irreducible over F_i). By the multiplicativity of extension degrees, we see that

$$[K : F] = [F_k : F_{k-1}][F_{k-1} : F_{k-2}] \cdots [F_1 : F_0]$$

is also finite, and $\leq n_1 n_2 \cdots n_k$.

This also gives a description of the elements of $F(\alpha_1, \alpha_2, \dots, \alpha_k)$. For simplicity, consider the case of the field $F(\alpha, \beta)$ where α and β are algebraic over F . Then the elements of this field are of the form

$$b_0 + b_1\beta + b_2\beta^2 + \cdots + b_{d-1}\beta^{d-1}$$

where $d = [F(\alpha)(\beta) : F(\alpha)]$ is the degree of β over $F(\alpha)$ (which may be strictly smaller than the degree of β over F), and where the coefficients b_0, b_1, \dots, b_{d-1} are elements of $F(\alpha)$. The coefficients $b_i \in F(\alpha)$, $i = 0, \dots, d-1$, are of the form

$$a_{0i} + a_{1i}\alpha + a_{2i}\alpha^2 + \cdots + a_{n-1i}\alpha^{n-1}$$

where $n = [F(\alpha) : F]$ is the degree of α over F and the a_{ij} are elements of F . Hence the elements of $F(\alpha, \beta)$ are of the form

$$\sum_{\substack{i=0,1,\dots,n-1 \\ j=0,1,\dots,d-1}} a_{ij}\alpha^i\beta^j \quad a_{ij} \in F.$$

Since $[F(\alpha, \beta) : F] = [F(\alpha, \beta) : F(\alpha)][F(\alpha) : F] = dn$, the elements $\alpha^i\beta^j$ are in fact an F basis for $F(\alpha, \beta)$.

In practice the field $F(\alpha)$ generated by the algebraic α is obtained by adjoining the element α to F and then “closing” the resulting set with respect to addition and multiplication, which amounts to adjoining the powers $\alpha^2, \alpha^3, \dots$ of α and taking linear combinations (with coefficients from F) of these elements. The process terminates when a power of α is a linear combination of lower powers of α which amounts to knowing the minimal polynomial for α . The previous discussion shows a similar process gives the field $F(\alpha, \beta)$ generated by two elements, and by recursion, the field generated by any finite number of algebraic elements. This shows in particular that “closing” with respect to addition and multiplication also closes with respect to division for algebraic elements (cf. Example 5 following Corollary 5 above). If the elements are not algebraic, one must also “close” with respect to inverses. The difficulty in this procedure is determining the degrees of the *relative* extensions — for example the degree d for $F(\alpha, \beta)$ over $F(\alpha)$ above, for which one has only an a priori upper bound (the degree of β over F).

This is the analogue of “closing” a set of elements in a group G to determine the subgroup they generate.

Examples

- (1) The extension $\mathbb{Q}(\sqrt[4]{2}, \sqrt{2})$ is simply the extension $\mathbb{Q}(\sqrt[4]{2})$ since $\sqrt{2}$ is already an element of this field. Put another way, the degree d of $\sqrt{2}$ over $\mathbb{Q}(\sqrt[4]{2})$ is 1, which is strictly smaller than the degree of $\sqrt{2}$ over \mathbb{Q} . We shall later have less obvious examples where this occurs.
- (2) Consider the field $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ generated over \mathbb{Q} by $\sqrt{2}$ and $\sqrt{3}$. Since $\sqrt{3}$ is of degree 2 over \mathbb{Q} the degree of the extension $\mathbb{Q}(\sqrt{2}, \sqrt{3})/\mathbb{Q}(\sqrt{2})$ is at most 2 and is precisely 2 if and only if $x^2 - 3$ is irreducible over $\mathbb{Q}(\sqrt{2})$. Since this polynomial is of degree 2, it is reducible only if it has a root, i.e., if and only if $\sqrt{3} \in \mathbb{Q}(\sqrt{2})$. Suppose $\sqrt{3} = a + b\sqrt{2}$ with $a, b \in \mathbb{Q}$. Squaring this we obtain $3 = (a^2 + 2b^2) + 2ab\sqrt{2}$. If $ab \neq 0$, then we can solve this equation for $\sqrt{2}$ in terms of a and b which implies that $\sqrt{2}$ is rational, which it is not. If $b = 0$, then we would have that $\sqrt{3} = a$ is rational, a contradiction. Finally, if $a = 0$, we have $\sqrt{3} = b\sqrt{2}$ and multiplying both sides by $\sqrt{2}$ we see that $\sqrt{6}$ would be rational, again a contradiction. This shows $\sqrt{3} \notin \mathbb{Q}(\sqrt{2})$, proving

$$[\mathbb{Q}(\sqrt{2}, \sqrt{3}) : \mathbb{Q}] = 4.$$

Elements in this field (by “closing” 1, $\sqrt{2}$, $\sqrt{3}$) include 1, $\sqrt{2}$, $\sqrt{3}$, $\sqrt{6}$ and by the computations above, these form a basis for this field:

$$\mathbb{Q}(\sqrt{2}, \sqrt{3}) = \{a + b\sqrt{2} + c\sqrt{3} + d\sqrt{6} \mid a, b, c, d \in \mathbb{Q}\}.$$

We can now characterize the finite extensions of a field F :

Theorem 17. The extension K/F is finite if and only if K is generated by a finite number of algebraic elements over F . More precisely, a field generated over F by a finite number of algebraic elements of degrees n_1, n_2, \dots, n_k is algebraic of degree $\leq n_1 n_2 \cdots n_k$.

Proof: If K/F is finite of degree n , let $\alpha_1, \alpha_2, \dots, \alpha_n$ be a basis for K as a vector space over F . By Corollary 15, $[F(\alpha_i) : F]$ divides $[K : F] = n$ for $i = 1, 2, \dots, n$, so

that Proposition 12 implies each α_i is algebraic over F . Since K is obviously generated over F by $\alpha_1, \alpha_2, \dots, \alpha_n$, we see that K is generated by a finite number of algebraic elements over F . The converse was proved above. The second statement of the theorem is immediate from Corollary 13 and the computation above.

The first example above shows that the inequality for the degree of the extension given in the theorem may be strict. We remark that information helpful in the determination of this degree can often be obtained by determining subfields and then applying Corollary 15.

Corollary 18. Suppose α and β are algebraic over F . Then $\alpha \pm \beta, \alpha\beta, \alpha/\beta$ (for $\beta \neq 0$), (in particular α^{-1} for $\alpha \neq 0$) are all algebraic.

Proof: All of these elements lie in the extension $F(\alpha, \beta)$, which is finite over F by the theorem, hence they are algebraic by Corollary 13.

Corollary 19. Let L/F be an arbitrary extension. Then the collection of elements of L that are algebraic over F form a subfield K of L .

Proof: This is immediate from the previous corollary.

Examples

- (1) Consider the extension \mathbb{C}/\mathbb{Q} and let $\overline{\mathbb{Q}}$ denote the subfield of all elements in \mathbb{C} that are algebraic over \mathbb{Q} . In particular, the elements $\sqrt[n]{2}$ (the positive n^{th} roots of 2 in \mathbb{R}) are all elements of $\overline{\mathbb{Q}}$, so that $[\overline{\mathbb{Q}} : \mathbb{Q}] \geq n$ for all integers $n > 1$. Hence $\overline{\mathbb{Q}}$ is an *infinite* algebraic extension of \mathbb{Q} , called the field of *algebraic numbers*.
- (2) Consider the field $\overline{\mathbb{Q}} \cap \mathbb{R}$, the subfield of \mathbb{R} consisting of elements algebraic over \mathbb{Q} . The field \mathbb{Q} is *countable*. The number of polynomials in $\mathbb{Q}[x]$ of any given degree n is therefore also countable (since such a polynomial is determined by specifying $n+1$ coefficients from \mathbb{Q}). Since these polynomials have at most n roots in \mathbb{R} , the number of algebraic elements of \mathbb{R} of degree n is countable. Finally, the collection of all algebraic elements in \mathbb{R} is the countable union (indexed by n) of countable sets, hence is countable. Since \mathbb{R} is uncountable, it follows that there exist (in fact many) elements of \mathbb{R} which are not algebraic, i.e., are transcendental, over \mathbb{Q} . In particular the subfield $\overline{\mathbb{Q}} \cap \mathbb{R}$ of algebraic elements of \mathbb{R} is a *proper* subfield of \mathbb{R} , so also $\overline{\mathbb{Q}}$ is a proper subfield of \mathbb{C} .

It is extremely difficult in general to prove that a given real number is not algebraic. For example, it is known (these are theorems) that $\pi = 3.14159\dots$ and $e = 2.71828\dots$ are transcendental elements of \mathbb{R} . Even the proofs that these elements are not *rational* are not too easy.

Theorem 20. If K is algebraic over F and L is algebraic over K , then L is algebraic over F .

Proof: Let α be any element of L . Then α is algebraic over K , so α satisfies some polynomial equation

$$a_n\alpha^n + a_{n-1}\alpha^{n-1} + \cdots + a_1\alpha + a_0 = 0$$

where the coefficients a_0, a_1, \dots, a_n are in K . Consider the field $F(\alpha, a_0, a_1, \dots, a_n)$ generated over F by α and the coefficients of this polynomial. Since K/F is algebraic, the elements a_0, a_1, \dots, a_n are algebraic over F , so the extension $F(a_0, a_1, \dots, a_n)/F$ is finite by Theorem 17. By the equation above, we see that α generates an extension of this field of degree at most n , since its minimal polynomial over this field is a divisor of the polynomial above. Therefore

$$[F(\alpha, a_0, a_1, \dots, a_n) : F] = [F(\alpha, a_0, \dots, a_n) : F(a_0, \dots, a_n)][F(a_0, \dots, a_n) : F]$$

is also finite and $F(\alpha, a_0, a_1, \dots, a_n)/F$ is an algebraic extension. In particular the element α is algebraic over F , which proves that L is algebraic over F .

The subfield $F(\alpha_1, \alpha_2, \dots, \alpha_k)$ generated by a finite set of elements $\alpha_1, \alpha_2, \dots, \alpha_k$ of a field K contains each of the fields $F(\alpha_i)$, $i = 1, 2, \dots, k$. By the definitions, it is also the smallest subfield of K containing these fields.

Definition. Let K_1 and K_2 be two subfields of a field K . Then the *composite field* of K_1 and K_2 , denoted $K_1 K_2$, is the smallest subfield of K containing both K_1 and K_2 . Similarly, the composite of any collection of subfields of K is the smallest subfield containing all the subfields.

Note that the composite $K_1 K_2$ can also be described as the intersection of all the subfields of K containing both K_1 and K_2 and similarly for the composite of more than two fields, analogous to the subgroup generated by a subset of a group (cf. Section 2.4).

Example

The composite of the two fields $\mathbb{Q}(\sqrt{2})$ and $\mathbb{Q}(\sqrt[3]{2})$ is the field $\mathbb{Q}(\sqrt[6]{2})$. This is because this field contains both of these subfields ($(\sqrt[6]{2})^3 = \sqrt{2}$ and $(\sqrt[6]{2})^2 = \sqrt[3]{2}$) and conversely, any field containing both $\sqrt{2}$ and $\sqrt[3]{2}$ contains their quotient, which is $\sqrt[6]{2}$.

Suppose now that K_1 and K_2 are finite extensions of F in K . Let $\alpha_1, \alpha_2, \dots, \alpha_n$ be an F -basis for K_1 and let $\beta_1, \beta_2, \dots, \beta_m$ be an F -basis for K_2 (so that $[K_1 : F] = n$ and $[K_2 : F] = m$). Then it is clear that these give generators for the composite $K_1 K_2$ over F :

$$K_1 K_2 = F(\alpha_1, \alpha_2, \dots, \alpha_n, \beta_1, \beta_2, \dots, \beta_m).$$

Since $\alpha_1, \alpha_2, \dots, \alpha_n$ is an F -basis for K_1 any power α_i^k of one of the α 's is a *linear* combination with coefficients in F of the α 's and a similar statement holds for the β 's. It follows that the collection of linear combinations

$$\sum_{\substack{i=1,2,\dots,n \\ j=1,2,\dots,m}} a_{ij} \alpha_i \beta_j$$

with coefficients in F is *closed* under multiplication and addition since in a product of two such elements any higher powers of the α 's and β 's can be replaced by linear expressions. Hence, the elements $\alpha_i \beta_j$ for $i = 1, 2, \dots, n$ and $j = 1, 2, \dots, m$ span the composite extension $K_1 K_2$ over F . In particular, $[K_1 K_2 : F] \leq mn$. We summarize this as:

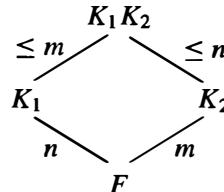
Proposition 21. Let K_1 and K_2 be two finite extensions of a field F contained in K . Then

$$[K_1 K_2 : F] \leq [K_1 : F][K_2 : F]$$

with equality if and only if an F -basis for one of the fields remains linearly independent over the other field. If $\alpha_1, \alpha_2, \dots, \alpha_n$ and $\beta_1, \beta_2, \dots, \beta_m$ are bases for K_1 and K_2 over F , respectively, then the elements $\alpha_i \beta_j$ for $i = 1, 2, \dots, n$ and $j = 1, 2, \dots, m$ span $K_1 K_2$ over F .

Proof: From $K_1 K_2 = F(\alpha_1, \alpha_2, \dots, \alpha_n, \beta_1, \beta_2, \dots, \beta_m) = K_1(\beta_1, \beta_2, \dots, \beta_m)$, we see as above that $\beta_1, \beta_2, \dots, \beta_m$ span $K_1 K_2$ over K_1 . Hence $[K_1 K_2 : K_1] \leq m = [K_2 : F]$ with equality if and only if these elements are linearly independent over K_1 . Since $[K_1 K_2 : F] = [K_1 K_2 : K_1][K_1 : F]$ this proves the proposition.

By the proposition (and its proof), we have the following diagram:



We shall have examples shortly where the inequality in the proposition is strict. One useful situation where one can be certain of equality is the following:

Corollary 22. Suppose that $[K_1 : F] = n$, $[K_2 : F] = m$ in Proposition 21, where n and m are relatively prime: $(n, m) = 1$. Then $[K_1 K_2 : F] = [K_1 : F][K_2 : F] = nm$.

Proof: In general the extension degree $[K_1 K_2 : F]$ is divisible by both n and m since K_1 and K_2 are subfields of $K_1 K_2$, hence is divisible by their least common multiple. In this case, since $(n, m) = 1$, this means $[K_1 K_2 : F]$ is divisible by nm , which together with the inequality $[K_1 K_2 : F] \leq nm$ of the proposition proves the corollary.

Example

The composite of the two fields $\mathbb{Q}(\sqrt{2})$ and $\mathbb{Q}(\sqrt[3]{2})$ is of degree 6 over \mathbb{Q} , which we determined earlier by actually computing the composite $\mathbb{Q}(\sqrt[6]{2})$.

EXERCISES

- Let \mathbb{F} be a finite field of characteristic p . Prove that $|\mathbb{F}| = p^n$ for some positive integer n .
- Let $g(x) = x^2 + x - 1$ and let $h(x) = x^3 - x + 1$. Obtain fields of 4, 8, 9 and 27 elements by adjoining a root of $f(x)$ to the field F where $f(x) = g(x)$ or $h(x)$ and $F = \mathbb{F}_2$ or \mathbb{F}_3 . Write down the multiplication tables for the fields with 4 and 9 elements and show that the nonzero elements form a cyclic group.
- Determine the minimal polynomial over \mathbb{Q} for the element $1 + i$.

4. Determine the degree over \mathbb{Q} of $2 + \sqrt{3}$ and of $1 + \sqrt[3]{2} + \sqrt[3]{4}$.
5. Let $F = \mathbb{Q}(i)$. Prove that $x^3 - 2$ and $x^3 - 3$ are irreducible over F .
6. Prove directly from the definitions that the field $F(\alpha_1, \alpha_2, \dots, \alpha_n)$ is the composite of the fields $F(\alpha_1), F(\alpha_2), \dots, F(\alpha_n)$.
7. Prove that $\mathbb{Q}(\sqrt{2} + \sqrt{3}) = \mathbb{Q}(\sqrt{2}, \sqrt{3})$ [one inclusion is obvious, for the other consider $(\sqrt{2} + \sqrt{3})^2$, etc.]. Conclude that $[\mathbb{Q}(\sqrt{2} + \sqrt{3}) : \mathbb{Q}] = 4$. Find an irreducible polynomial satisfied by $\sqrt{2} + \sqrt{3}$.
8. Let F be a field of characteristic $\neq 2$. Let D_1 and D_2 be elements of F , neither of which is a square in F . Prove that $F(\sqrt{D_1}, \sqrt{D_2})$ is of degree 4 over F if $D_1 D_2$ is not a square in F and is of degree 2 over F otherwise. When $F(\sqrt{D_1}, \sqrt{D_2})$ is of degree 4 over F the field is called a *biquadratic extension of F* .
9. Let F be a field of characteristic $\neq 2$. Let a, b be elements of the field F with b not a square in F . Prove that a necessary and sufficient condition for $\sqrt{a + \sqrt{b}} = \sqrt{m} + \sqrt{n}$ for some m and n in F is that $a^2 - b$ is a square in F . Use this to determine when the field $\mathbb{Q}(\sqrt{a + \sqrt{b}})$ ($a, b \in \mathbb{Q}$) is biquadratic over \mathbb{Q} .
10. Determine the degree of the extension $\mathbb{Q}(\sqrt{3 + 2\sqrt{2}})$ over \mathbb{Q} .
11. (a) Let $\sqrt{3+4i}$ denote the square root of the complex number $3 + 4i$ that lies in the first quadrant and let $\sqrt{3-4i}$ denote the square root of $3 - 4i$ that lies in the fourth quadrant. Prove that $[\mathbb{Q}(\sqrt{3+4i} + \sqrt{3-4i}) : \mathbb{Q}] = 1$.
(b) Determine the degree of the extension $\mathbb{Q}(\sqrt{1+\sqrt{-3}} + \sqrt{1-\sqrt{-3}})$ over \mathbb{Q} .
12. Suppose the degree of the extension K/F is a prime p . Show that any subfield E of K containing F is either K or F .
13. Suppose $F = \mathbb{Q}(\alpha_1, \alpha_2, \dots, \alpha_n)$ where $\alpha_i^2 \in \mathbb{Q}$ for $i = 1, 2, \dots, n$. Prove that $\sqrt[3]{2} \notin F$.
14. Prove that if $[F(\alpha) : F]$ is odd then $F(\alpha) = F(\alpha^2)$.
15. A field F is said to be *formally real* if -1 is not expressible as a sum of squares in F . Let F be a formally real field, let $f(x) \in F[x]$ be an irreducible polynomial of odd degree and let α be a root of $f(x)$. Prove that $F(\alpha)$ is also formally real. [Pick α a counterexample of minimal degree. Show that $-1 + f(x)g(x) = (p_1(x))^2 + \dots + (p_m(x))^2$ for some $p_i(x), g(x) \in F[x]$ where $g(x)$ has odd degree $< \deg f$. Show that some root β of g has odd degree over F and $F(\beta)$ is not formally real, violating the minimality of α .]
16. Let K/F be an algebraic extension and let R be a ring contained in K and containing F . Show that R is a subfield of K containing F .
17. Let $f(x)$ be an irreducible polynomial of degree n over a field F . Let $g(x)$ be any polynomial in $F[x]$. Prove that every irreducible factor of the composite polynomial $f(g(x))$ has degree divisible by n .
18. Let k be a field and let $k(x)$ be the field of rational functions in x with coefficients from k . Let $t \in k(x)$ be the rational function $\frac{P(x)}{Q(x)}$ with relatively prime polynomials $P(x), Q(x) \in k[x]$, with $Q(x) \neq 0$. Then $k(x)$ is an extension of $k(t)$ and to compute its degree it is necessary to compute the minimal polynomial with coefficients in $k(t)$ satisfied by x .
 - (a) Show that the polynomial $P(X) - tQ(X)$ in the variable X and coefficients in $k(t)$ is irreducible over $k(t)$ and has x as a root. [By Gauss' Lemma this polynomial is irreducible in $(k(t))[X]$ if and only if it is irreducible in $(k[t])[X]$. Then note that $(k[t])[X] = (k[X])[t]$.]

- (b) Show that the degree of $P(X) - tQ(X)$ as a polynomial in X with coefficients in $k(t)$ is the maximum of the degrees of $P(x)$ and $Q(x)$.
- (c) Show that $[k(x) : k(t)] = [k(x) : k(\frac{P(x)}{Q(x)})] = \max(\deg P(x), \deg Q(x))$.
19. Let K be an extension of F of degree n .
- For any $\alpha \in K$ prove that α acting by left multiplication on K is an F -linear transformation of K .
 - Prove that K is isomorphic to a subfield of the ring of $n \times n$ matrices over F , so the ring of $n \times n$ matrices over F contains an isomorphic copy of every extension of F of degree $\leq n$.
20. Show that if the matrix of the linear transformation “multiplication by α ” considered in the previous exercise is A then α is a root of the characteristic polynomial for A . This gives an effective procedure for determining an equation of degree n satisfied by an element α in an extension of F of degree n . Use this procedure to obtain the monic polynomial of degree 3 satisfied by $\sqrt[3]{2}$ and by $1 + \sqrt[3]{2} + \sqrt[3]{4}$.
21. Let $K = \mathbb{Q}(\sqrt{D})$ for some squarefree integer D . Let $\alpha = a + b\sqrt{D}$ be an element of K . Use the basis $1, \sqrt{D}$ for K as a vector space over \mathbb{Q} and show that the matrix of the linear transformation “multiplication by α ” on K considered in the previous exercises has the matrix $\begin{pmatrix} a & bD \\ b & a \end{pmatrix}$. Prove directly that the map $a + b\sqrt{D} \mapsto \begin{pmatrix} a & bD \\ b & a \end{pmatrix}$ is an isomorphism of the field K with a subfield of the ring of 2×2 matrices with coefficients in \mathbb{Q} .
22. Let K_1 and K_2 be two finite extensions of a field F contained in the field K . Prove that the F -algebra $K_1 \otimes_F K_2$ is a field if and only if $[K_1 K_2 : F] = [K_1 : F][K_2 : F]$.

13.3 CLASSICAL STRAIGHTEDGE AND COMPASS CONSTRUCTIONS

As a simple application of the results we have obtained on algebraic extensions, and in particular on the multiplicativity of extension degrees, we can answer (in the negative) the following geometric problems posed by the Greeks:

- I. (*Doubling the Cube*) Is it possible using only straightedge and compass to construct a cube with precisely twice the volume of a given cube?
- II. (*Trisecting an Angle*) Is it possible using only straightedge and compass to trisect any given angle θ ?
- III. (*Squaring the Circle*) Is it possible using only straightedge and compass to construct a square whose area is precisely the area of a given circle?

To answer these questions we must translate the construction of lengths by compass and straightedge into algebraic terms. Let 1 denote a fixed given unit distance. Then any distance is determined by its length $a \in \mathbb{R}$, which allows us to view geometric distances as elements of the real numbers \mathbb{R} . Using the given unit distance 1 to define the scale on the axes, we can then construct the usual Cartesian plane \mathbb{R}^2 and view all of our constructions as occurring in \mathbb{R}^2 . A point $(x, y) \in \mathbb{R}^2$ is then constructible starting with the given distance 1 if and only if its coordinates x and y are constructible elements of \mathbb{R} . The problems above then amount to determining whether particular lengths in \mathbb{R} can be obtained by compass and straightedge constructions from a fixed

unit distance. The collection of such real numbers together with their negatives will be called the *constructible* elements of \mathbb{R} , and we shall not distinguish between the lengths that are constructible and the real numbers that are constructible.

Each straightedge and compass construction consists of a series of operations of the following four types: (1) connecting two given points by a straight line, (2) finding a point of intersection of two straight lines, (3) drawing a circle with given radius and center, and (4) finding the point(s) of intersection of a straight line and a circle or the intersection of two circles.

It is an elementary fact from geometry that if two lengths a and b are given one may construct using straightedge and compass the lengths $a \pm b$, ab and a/b (the first two are clear and the latter two are given by the construction of parallel lines (Figure 1)).

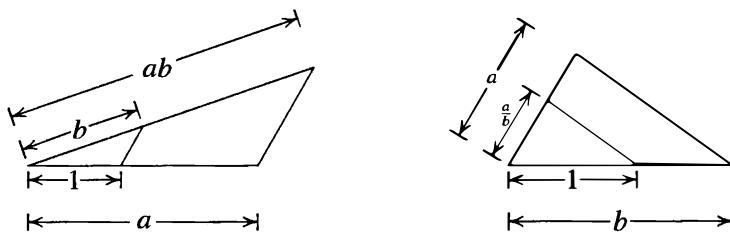


Fig. 1

It is also an elementary geometry construction to construct \sqrt{a} if a is given: construct the circle with diameter $1 + a$ and erect the perpendicular to the diameter as indicated in Figure 2. Then \sqrt{a} is the length of this perpendicular.

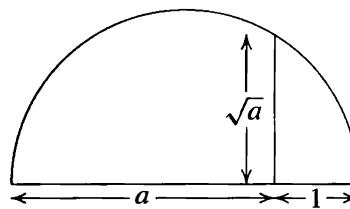


Fig. 2

It follows that straightedge and compass constructions give all the algebraic operations of addition, subtraction, multiplication and division (by nonzero elements) in the reals so the collection of constructible elements is a *subfield* of \mathbb{R} . One can also take square roots of constructible elements. We shall now see that these are essentially the only operations possible.

From the given length 1 it is possible to construct by these operations all the rational numbers \mathbb{Q} . Hence we may construct all of the points $(x, y) \in \mathbb{R}^2$ whose coordinates are rational. We may construct additional elements of \mathbb{R} by taking square roots, so the collection of elements constructible from 1 of \mathbb{R} form a field strictly larger than \mathbb{Q} .

The usual formula (“two point form”) for the straight line connecting two points with coordinates in some field F gives an equation for the line of the form $ax + by - c = 0$ with $a, b, c \in F$. Solving two such equations simultaneously to determine the point of intersection of two such lines gives solutions also in F . It follows that if the coordinates

of two points lie in the field F then straightedge constructions alone will not produce additional points whose coordinates are not also in F .

A compass construction (type (3) or (4) above) defines points obtained by the intersection of a circle with either a straight line or another circle. A circle with center (h, k) and radius r has equation

$$(x - h)^2 + (y - k)^2 = r^2$$

so when we consider the effect of compass constructions on elements of a field F we are considering simultaneous solutions of such an equation with a linear equation $ax + by - c = 0$ where $a, b, c, h, k, r \in F$, or the simultaneous solutions of two quadratic equations.

In the case of a linear equation and the equation for the circle, solving for y , say, in the linear equation and substituting gives a *quadratic* equation for x (and y is given linearly in terms of x). Hence the coordinates of the point of intersection are at worst in a *quadratic extension* of F .

In the case of the intersection of two circles, say

$$(x - h)^2 + (y - k)^2 = r^2$$

$$\text{and } (x - h')^2 + (y - k')^2 = r'^2,$$

subtraction of the second equation from the first shows that we have the same intersection by considering the two equations

$$(x - h)^2 + (y - k)^2 = r^2$$

$$\text{and } 2(h' - h)x + 2(k' - k)y = r^2 - h^2 - k^2 - r'^2 + h'^2 + k'^2$$

which is the intersection of a circle and a straight line (the straight line connecting the two points of intersection, in fact) of the type just considered.

It follows that if a collection of constructible elements is given, then one can construct all the elements in the subfield F of \mathbb{R} generated by these elements and that any straightedge and compass operation on elements of F produces elements in at worst a *quadratic* extension of F . Since quadratic extensions have degree 2 and extension degrees are multiplicative, it follows that if $\alpha \in \mathbb{R}$ is obtained from elements in a field F by a (finite) series of straightedge and compass operations then α is an element of an extension K of F of degree a power of 2: $[K : F] = 2^m$ for some m . Since $[F(\alpha) : F]$ divides this extension degree, it must also be a power of 2.

Proposition 23. If the element $\alpha \in \mathbb{R}$ is obtained from a field $F \subset \mathbb{R}$ by a series of compass and straightedge constructions then $[F(\alpha) : F] = 2^k$ for some integer $k \geq 0$.

Theorem 24. None of the classical Greek problems: (I) Doubling the Cube, (II) Trisecting an Angle, and (III) Squaring the Circle, is possible.

Proof: (I) Doubling the cube amounts to constructing $\sqrt[3]{2}$ in the reals starting with the unit 1. Since $[\mathbb{Q}(\sqrt[3]{2}) : \mathbb{Q}] = 3$ is not a power of 2, this is impossible.

(II) If an angle θ can be constructed, then determining the point at distance 1 from the origin and angle θ from the positive x axis in \mathbb{R}^2 shows that $\cos \theta$ (the x -coordinate

of this point) can be constructed (so then $\sin \theta$ can also be constructed). Conversely if $\cos \theta$, then $\sin \theta$, can be constructed, the point with those coordinates gives the angle θ .

The problem of trisecting the angle θ is then equivalent to the problem: given $\cos \theta$ construct $\cos \theta/3$.

To see that this is not always possible (it is certainly occasionally possible, for example for $\theta = 180^\circ$), consider $\theta = 60^\circ$. Then $\cos \theta = \frac{1}{2}$. By the triple angle formula for cosines:

$$\cos \theta = 4\cos^3 \theta/3 - 3\cos \theta/3,$$

substituting $\theta = 60^\circ$, we see that $\beta = \cos 20^\circ$ satisfies the equation

$$4\beta^3 - 3\beta - 1/2 = 0$$

or $8(\beta)^3 - 6\beta - 1 = 0$. This can be written $(2\beta)^3 - 3(2\beta) - 1 = 0$. Let $\alpha = 2\beta$. Then α is a real number between 0 and 2 satisfying the equation

$$\alpha^3 - 3\alpha - 1 = 0.$$

But we considered this equation in the last section and determined $[\mathbb{Q}(\alpha) : \mathbb{Q}] = 3$, and as before we see that α is not constructible.

(III) Squaring the circle is equivalent to determining whether the real number $\pi = 3.14159 \dots$ is constructible. As mentioned previously, it is a difficult problem even to prove that this number is not rational. It is in fact transcendental (which we shall assume without proof), so that $[\mathbb{Q}(\pi) : \mathbb{Q}]$ is not even finite, much less a power of 2, showing the impossibility of squaring the circle by straightedge and compass.

Remark: The proof above shows that $\cos 20^\circ$ and $\sin 20^\circ$ cannot be constructed. The question arises as to which integer angles (measured in degrees) are constructible? The angles 1° and 2° are not constructible, since otherwise the addition formulae for sines and cosines would give the constructibility for 20° . On the other hand, elementary geometric constructions (of the regular 5-gon for an angle of 72° and the equilateral triangle for an angle of 60°) together with the addition formulae and the half-angle formulae show that $\cos 3^\circ$ and $\sin 3^\circ$ are constructible. It follows from this that the trigonometric functions of an integer degree angle are constructible precisely when the angle is a multiple of 3° . Explicitly,

$$\begin{aligned}\cos 3^\circ &= \frac{1}{8}(\sqrt{3} + 1)\sqrt{5 + \sqrt{5}} + \frac{1}{16}(\sqrt{6} - \sqrt{2})(\sqrt{5} - 1) \\ \sin 3^\circ &= \frac{1}{16}(\sqrt{6} + \sqrt{2})(\sqrt{5} - 1) - \frac{1}{8}(\sqrt{3} - 1)\sqrt{5 + \sqrt{5}},\end{aligned}$$

showing that these are obtained from \mathbb{Q} by successive extractions of square roots and field operations.

After discussing the cyclotomic fields in Section 14.5 we shall consider another classical geometric question: “which regular n -gons can be constructed by straightedge and compass?” (cf. Proposition 14.29).

We have been careful here to consider constructions using a *straightedge* rather than a *ruler*, the distinction being that a ruler has marks on it. If one uses a ruler, it is

possible to construct many additional algebraic elements. For example, suppose θ is a given angle and the unit distance 1 is marked on the ruler. Draw a circle of radius 1 with central angle θ as shown in Figure 3 and then slide the ruler until the distance between points A and B on the circle is 1. Then some elementary geometry shows that (cf. the exercises) the angle α indicated is $\theta/3$, i.e., this construction (due to Archimedes) trisects θ . In particular, the second classical problem in Theorem 24 (Trisecting an Angle) can be solved with *ruler* and compass.

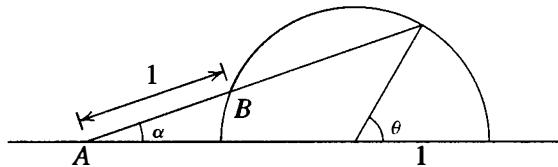


Fig. 3

The first of the classical problems in Theorem 24 (Duplication of the Cube), which amounts to the construction of $\sqrt[3]{2}$, can also be solved with ruler and compass. The following gives a construction for $k^{1/3}$ for any given positive real k which is less than 1. This construction was shown to us by J.H. Conway.

Drawing a circle of radius 1 and using the point $A = (k, 0)$ as center, construct the point $B = (0, \sqrt{1-k^2})$. Dividing this distance by 3, construct the point $(0, -\frac{1}{3}\sqrt{1-k^2})$ and draw the line connecting this point with A . Slide the ruler with marked unit length 1 so that it passes through the point B and so that the distance from the intersection point C to the intersection point D with the x -axis is of length 1, as indicated in Figure 4.

Then the distance between A and D is $2k^{1/3}$ and the distance between B and C is $2k^{2/3}$ (cf. the exercises).

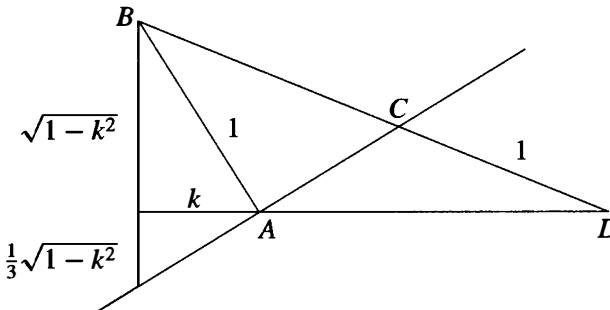


Fig. 4

EXERCISES

1. Prove that it is impossible to construct the regular 9-gon.
2. Prove that Archimedes' construction actually trisects the angle θ . [Note the isosceles triangles in Figure 5 to prove that $\beta = \gamma = 2\alpha$.]

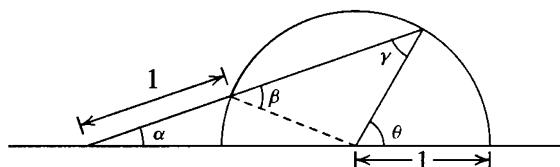


Fig. 5

3. Prove that Conway's construction indicated in the text actually constructs $2k^{1/3}$ and $2k^{2/3}$. [One method: let (x, y) be the coordinates of the point C , a the distance from B to C and b the distance from A to D ; use similar triangles to prove (a) $\frac{y}{1} = \frac{\sqrt{1-k^2}}{1+a}$, (b) $\frac{x}{a} = \frac{b+k}{1+a}$, (c) $\frac{y}{x-k} = \frac{\sqrt{1-k^2}}{3k}$, and also show that (d) $(1-k^2)+(b+k)^2 = (1+a)^2$; solve these equations for a and b .]
4. The construction of the regular 7-gon amounts to the constructibility of $\cos(2\pi/7)$. We shall see later (Section 14.5 and Exercise 2 of Section 14.7) that $\alpha = 2\cos(2\pi/7)$ satisfies the equation $x^3 + x^2 - 2x - 1 = 0$. Use this to prove that the regular 7-gon is not constructible by straightedge and compass.
5. Use the fact that $\alpha = 2\cos(2\pi/5)$ satisfies the equation $x^2 + x - 1 = 0$ to conclude that the regular 5-gon is constructible by straightedge and compass.

13.4 SPLITTING FIELDS AND ALGEBRAIC CLOSURES

Let F be a field.

If $f(x)$ is any polynomial in $F[x]$ then we have seen in Section 2 that there exists a field K which can (by identifying F with an isomorphic copy of F) be considered an extension of F in which $f(x)$ has a root α . This is equivalent to the statement that $f(x)$ has a linear factor $x - \alpha$ in $K[x]$ (this is Proposition 9 of Chapter 9).

Definition. The extension field K of F is called a *splitting field* for the polynomial $f(x) \in F[x]$ if $f(x)$ factors completely into linear factors (or *splits completely*) in $K[x]$ and $f(x)$ does not factor completely into linear factors over any proper subfield of K containing F .

If $f(x)$ is of degree n , then $f(x)$ has at most n roots in F (Proposition 17 of Chapter 9) and has precisely n roots (counting multiplicities) in F if and only if $f(x)$ splits completely in $F[x]$.

Theorem 25. For any field F , if $f(x) \in F[x]$ then there exists an extension K of F which is a splitting field for $f(x)$.

Proof: We first show that there is an extension E of F over which $f(x)$ splits completely into linear factors by induction on the degree n of $f(x)$. If $n = 1$, then take $E = F$. Suppose now that $n > 1$. If the irreducible factors of $f(x)$ over F are all of degree 1, then F is the splitting field for $f(x)$ and we may take $E = F$. Otherwise, at least one of the irreducible factors, say $p(x)$ of $f(x)$ in $F[x]$ is of degree at least 2. By Theorem 3 there is an extension E_1 of F containing a root α of $p(x)$. Over E_1 the polynomial $f(x)$ has the linear factor $x - \alpha$. The degree of the remaining factor $f_1(x)$ of $f(x)$ is $n - 1$, so by induction there is an extension E of E_1 containing all the roots of $f_1(x)$. Since $\alpha \in E$, E is an extension of F containing all the roots of $f(x)$. Now let K be the intersection of all the subfields of E containing F which also contain all the roots of $f(x)$. Then K is a field which is a splitting field for $f(x)$.

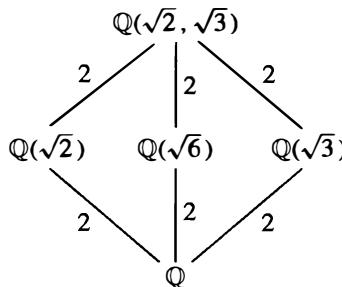
We shall see shortly that any two splitting fields for $f(x)$ are isomorphic (which extends Theorem 8), so (by abuse) we frequently refer to *the* splitting field of a polynomial.

Definition. If K is an algebraic extension of F which is the splitting field over F for a collection of polynomials $f(x) \in F[x]$ then K is called a *normal* extension of F .

We shall generally use the term “splitting field” rather than “normal extension” (cf. also Section 14.9).

Examples

- (1) The splitting field for $x^2 - 2$ over \mathbb{Q} is just $\mathbb{Q}(\sqrt{2})$, since the two roots are $\pm\sqrt{2}$ and $-\sqrt{2} \in \mathbb{Q}(\sqrt{2})$.
- (2) The splitting field for $(x^2 - 2)(x^2 - 3)$ is the field $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ generated over \mathbb{Q} by $\sqrt{2}$ and $\sqrt{3}$ since the roots of the polynomial are $\pm\sqrt{2}, \pm\sqrt{3}$. We have already seen that this is an extension of degree 4 over \mathbb{Q} and we have the following diagram of known subfields:



- (3) The splitting field of $x^3 - 2$ over \mathbb{Q} is not just $\mathbb{Q}(\sqrt[3]{2})$ since as previously noted the three roots of this polynomial in \mathbb{C} are

$$\sqrt[3]{2}, \quad \sqrt[3]{2} \left(\frac{-1 + i\sqrt{3}}{2} \right), \quad \sqrt[3]{2} \left(\frac{-1 - i\sqrt{3}}{2} \right)$$

and the latter two roots are not elements of $\mathbb{Q}(\sqrt[3]{2})$, since the elements of this field are of the form $a + b\sqrt[3]{2} + c\sqrt[3]{4}$ with rational a, b, c and all such numbers are real.

The splitting field K of this polynomial is obtained by adjoining all three of these roots to \mathbb{Q} . Note that since K contains the first two roots above, then it contains their quotient $\frac{-1 + \sqrt{-3}}{2}$ hence K contains the element $\sqrt{-3}$. On the other hand, any field containing $\sqrt[3]{2}$ and $\sqrt{-3}$ contains all three of the roots above. It follows that

$$K = \mathbb{Q}(\sqrt[3]{2}, \sqrt{-3})$$

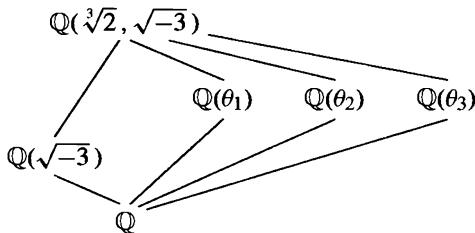
is the splitting field of $x^3 - 2$ over \mathbb{Q} . Since $\sqrt{-3}$ satisfies the equation $x^2 + 3 = 0$, the degree of this extension over $\mathbb{Q}(\sqrt[3]{2})$ is at most 2, hence must be 2 since we observed above that $\mathbb{Q}(\sqrt[3]{2})$ is not the splitting field. It follows that

$$[\mathbb{Q}(\sqrt[3]{2}, \sqrt{-3}) : \mathbb{Q}] = 6.$$

Note that we could have proceeded slightly differently at the end by noting that $\mathbb{Q}(\sqrt{-3})$ is a subfield of K , so that the index $[\mathbb{Q}(\sqrt{-3}) : \mathbb{Q}] = 2$ divides $[K : \mathbb{Q}]$.

Since this extension degree is also divisible by 3 (because $\mathbb{Q}(\sqrt[3]{2}) \subset K$), the degree is divisible by 6, hence must be 6.

This gives us the diagram of known subfields:



where

$$\theta_1 = \sqrt[3]{2}, \quad \theta_2 = \sqrt[3]{2} \left(\frac{-1 + i\sqrt{3}}{2} \right), \quad \theta_3 = \sqrt[3]{2} \left(\frac{-1 - i\sqrt{3}}{2} \right).$$

- (4) One must be careful in computing splitting fields. The splitting field for the polynomial $x^4 + 4$ over \mathbb{Q} is smaller than one might at first suspect. In fact this polynomial factors over \mathbb{Q} :

$$\begin{aligned} x^4 + 4 &= x^4 + 4x^2 + 4 - 4x^2 = (x^2 + 2)^2 - 4x^2 \\ &= (x^2 + 2x + 2)(x^2 - 2x + 2) \end{aligned}$$

where these two factors are irreducible (Eisenstein again). Solving for the roots of the two factors by the quadratic formula, we find the four roots

$$\pm 1 \pm i$$

so that the splitting field of this polynomial is just the field $\mathbb{Q}(i)$, an extension of degree 2 of \mathbb{Q} .

In general, if $f(x) \in F[x]$ is a polynomial of degree n , then adjoining one root of $f(x)$ to F generates an extension F_1 of degree at most n (and equal to n if and only if $f(x)$ is irreducible). Over F_1 the polynomial $f(x)$ now has at least one linear factor, so that any other root of $f(x)$ satisfies an equation of degree at most $n - 1$ over F_1 . Adjoining such a root to F_1 we therefore obtain an extension of degree at most $n - 1$ of F_1 , etc. Using the multiplicativity of extension degrees, this proves

Proposition 26. A splitting field of a polynomial of degree n over F is of degree at most $n!$ over F .

As the examples above show, the degree of a splitting field may be smaller than $n!$. It will be proved later using Galois Theory that a “general” polynomial of degree n (in a well defined sense) over \mathbb{Q} has a splitting field of degree $n!$, so this may be viewed as the “generic” situation (although most of the interesting examples we shall consider have splitting fields of smaller degree).

Example: (Splitting Field of $x^n - 1$: Cyclotomic Fields)

Consider the splitting field of the polynomial $x^n - 1$ over \mathbb{Q} . The roots of this polynomial are called the n^{th} roots of unity.

Recall that every nonzero complex number $a + bi \in \mathbb{C}$ can be written uniquely in the form

$$re^{i\theta} = r(\cos \theta + i \sin \theta) \quad r > 0, \quad 0 \leq \theta < 2\pi$$

which is simply representing the point $a + bi$ in the complex plane in terms of polar coordinates: r is the distance of (a, b) from the origin and θ is the angle made with the real positive axis.

Over \mathbb{C} there are n distinct solutions of the equation $x^n = 1$, namely the elements

$$e^{2\pi ki/n} = \cos\left(\frac{2\pi k}{n}\right) + i \sin\left(\frac{2\pi k}{n}\right)$$

for $k = 0, 1, \dots, n - 1$. These points are given geometrically by n equally spaced points starting with the point $(1, 0)$ (corresponding to $k = 0$) on a circle of radius 1 in the complex plane (see Figure 6). The fact that these are all n^{th} roots of unity is immediate, since

$$(e^{2\pi ki/n})^n = e^{(2\pi ki/n)n} = e^{2\pi ki} = 1.$$

It follows that \mathbb{C} contains a splitting field for $x^n - 1$ and we shall frequently view the splitting field for $x^n - 1$ over \mathbb{Q} as the field generated over \mathbb{Q} in \mathbb{C} by the numbers above.

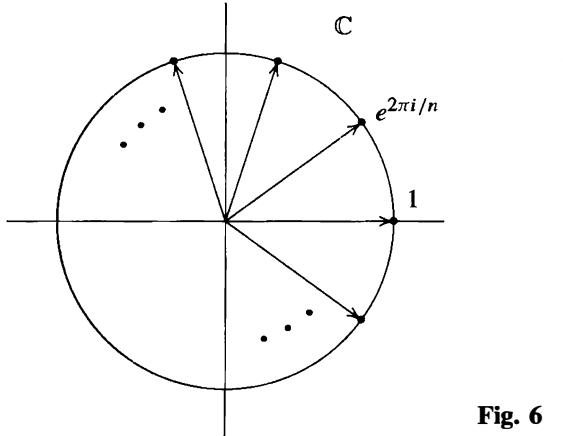


Fig. 6

In any abstract splitting field K/\mathbb{Q} for $x^n - 1$ the collection of n^{th} roots of unity form a group under multiplication since if $\alpha^n = 1$ and $\beta^n = 1$ then $(\alpha\beta)^n = 1$, so this subset of K^\times is closed under multiplication. It follows that this is a cyclic group (Proposition 18 of Chapter 9); we shall see that there are n distinct roots in K so it has order n .

Definition. A generator of the cyclic group of all n^{th} roots of unity is called a *primitive n^{th} root of unity*.

Let ζ_n denote a primitive n^{th} root of unity. The other primitive n^{th} roots of unity are then the elements ζ_n^a where $1 \leq a < n$ is an integer relatively prime to n , since these are the other generators for a cyclic group of order n . In particular there are precisely $\varphi(n)$ primitive n^{th} roots of unity, where $\varphi(n)$ denotes the Euler φ -function.

Over \mathbb{C} we can see all of this directly by letting

$$\zeta_n = e^{2\pi i/n}$$

(the first n^{th} root of unity counterclockwise from 1). Then all the other roots of unity are powers of ζ_n :

$$e^{2\pi ki/n} = \zeta_n^k$$

so that ζ_n is one possible generator for the multiplicative group of n^{th} roots of unity. When we view the roots of unity in \mathbb{C} we shall usually use ζ_n to denote this choice of a primitive n^{th} root of unity. The primitive roots of unity in \mathbb{C} for some small values of n are

$$\zeta_1 = 1$$

$$\zeta_2 = -1$$

$$\zeta_3 = \frac{-1 + i\sqrt{3}}{2}$$

$$\zeta_4 = i$$

$$\zeta_5 = \frac{\sqrt{5} - 1}{4} + i\left(\frac{\sqrt{10 + 2\sqrt{5}}}{4}\right)$$

$$\zeta_6 = \frac{1 + i\sqrt{3}}{2}$$

$$\zeta_8 = \frac{\sqrt{2}}{2} + i\frac{\sqrt{2}}{2}$$

(these formulas follow from the elementary geometry of n -gons and in any case can be verified directly by raising them to the appropriate power).

The splitting field of $x^n - 1$ over \mathbb{Q} is the field $\mathbb{Q}(\zeta_n)$ and this field is given a name:

Definition. The field $\mathbb{Q}(\zeta_n)$ is called the *cyclotomic field of n^{th} roots of unity*.

Determining the degree of this extension requires some analysis of the minimal polynomial of ζ_n over \mathbb{Q} and will be postponed until later (Section 6). One important special case which we have in fact already considered is when $n = p$ is a *prime*. In this case, we have the factorization

$$x^p - 1 = (x - 1)(x^{p-1} + x^{p-2} + \cdots + x + 1)$$

and since $\zeta_p \neq 1$ it follows that ζ_p is a root of the polynomial

$$\Phi_p(x) = \frac{x^p - 1}{x - 1} = x^{p-1} + x^{p-2} + \cdots + x + 1$$

which we showed was irreducible in Section 9.4. It follows that $\Phi_p(x)$ is the minimal polynomial of ζ_p over \mathbb{Q} , so that

$$[\mathbb{Q}(\zeta_p) : \mathbb{Q}] = p - 1.$$

We shall see later that in general $[\mathbb{Q}(\zeta_n) : \mathbb{Q}] = \varphi(n)$, where $\varphi(n)$ is the Euler phi-function of n (so that $\varphi(p) = p - 1$).

Example: (Splitting Field of $x^p - 2$, p a prime)

Let p be a prime and consider the splitting field of $x^p - 2$. If α is a root of this equation, i.e., $\alpha^p = 2$, then $(\zeta\alpha)^p = 2$ where ζ is any p^{th} root of unity. Hence the solutions of this equation are

$$\zeta \sqrt[p]{2}, \quad \zeta \text{ a } p^{\text{th}} \text{ root of unity}$$

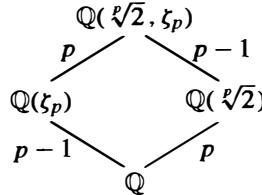
where as usual the symbol $\sqrt[p]{2}$ denotes the positive real p^{th} root of 2 if we wish to view these elements as complex numbers, and denotes any one solution of $x^p = 2$ if we view these roots abstractly. Since the ratio of the two solutions $\zeta_p \sqrt[p]{2}$ and $\sqrt[p]{2}$ for ζ_p , a primitive p^{th} root of unity is just ζ_p , the splitting field of $x^p - 2$ over \mathbb{Q} contains $\mathbb{Q}(\sqrt[p]{2}, \zeta_p)$. On the other hand, all the roots above lie in this field, so that the splitting field is precisely

$$\mathbb{Q}(\sqrt[p]{2}, \zeta_p).$$

This field contains the cyclotomic field of p^{th} roots of unity and is generated over it by $\sqrt[p]{2}$, hence is an extension of degree at most p . It follows that the degree of this extension over \mathbb{Q} is $\leq p(p-1)$. Since both $\mathbb{Q}(\sqrt[p]{2})$ and $\mathbb{Q}(\zeta_p)$ are subfields, the degree of the extension over \mathbb{Q} is divisible by p and by $p-1$. Since these two numbers are relatively prime it follows that the extension degree is divisible by $p(p-1)$ so that we must have

$$[\mathbb{Q}(\sqrt[p]{2}, \zeta_p) : \mathbb{Q}] = p(p-1)$$

(this is Corollary 22). Note in particular that we have proved $x^p - 2$ remains irreducible over $\mathbb{Q}(\zeta_p)$, which is not at all obvious. We have the following diagram of known subfields:



The special case $p = 3$ was Example 3 above, where we simply indicated the 3rd roots of unity explicitly.

We now return to the problem of proving it makes no difference how the splitting field of a polynomial $f(x)$ over a field F is constructed. As in Theorem 8 it is convenient to state the result for an arbitrary isomorphism $φ : F \xrightarrow{\sim} F'$ between two fields.

Theorem 27. Let $φ : F \xrightarrow{\sim} F'$ be an isomorphism of fields. Let $f(x) ∈ F[x]$ be a polynomial and let $f'(x) ∈ F'[x]$ be the polynomial obtained by applying $φ$ to the coefficients of $f(x)$. Let E be a splitting field for $f(x)$ over F and let E' be a splitting field for $f'(x)$ over F' . Then the isomorphism $φ$ extends to an isomorphism $σ : E \xrightarrow{\sim} E'$, i.e., $σ$ restricted to F is the isomorphism $φ$:

$$\begin{array}{ccc} σ : & E & \xrightarrow{\sim} E' \\ & | & | \\ φ : & F & \xrightarrow{\sim} F' \end{array}$$

Proof: We shall proceed by induction on the degree n of $f(x)$. As in the discussion before Theorem 8, recall that an isomorphism $φ$ from one field F to another field

F' induces a natural isomorphism between the polynomial rings $F[x]$ and $F'[x]$. In particular, if $f(x)$ and $f'(x)$ correspond to one another under this isomorphism then the irreducible factors of $f(x)$ in $F[x]$ correspond to the irreducible factors of $f'(x)$ in $F'[x]$.

If $f(x)$ has all its roots in F then $f(x)$ splits completely in $F[x]$ and $f'(x)$ splits completely in $F'[x]$ (with its linear factors being the images of the linear factors for $f(x)$). Hence $E = F$ and $E' = F'$, and in this case we may take $\sigma = \varphi$. This shows the result is true for $n = 1$ and in the case where all the irreducible factors of $f(x)$ have degree 1.

Assume now by induction that the theorem has been proved for any field F , isomorphism φ , and polynomial $f(x) \in F[x]$ of degree $< n$. Let $p(x)$ be an irreducible factor of $f(x)$ in $F[x]$ of degree at least 2 and let $p'(x)$ be the corresponding irreducible factor of $f'(x)$ in $F'[x]$. If $\alpha \in E$ is a root of $p(x)$ and $\beta \in E'$ is a root of $p'(x)$, then by Theorem 8 we can extend φ to an isomorphism $\sigma' : F(\alpha) \xrightarrow{\sim} F'(\beta)$:

$$\begin{array}{ccc} \sigma' : & F(\alpha) & \xrightarrow{\sim} F'(\beta) \\ & | & | \\ \varphi : & F & \xrightarrow{\sim} F'. \end{array}$$

Let $F_1 = F(\alpha)$, $F'_1 = F'(\beta)$, so that we have the isomorphism $\sigma' : F_1 \xrightarrow{\sim} F'_1$. We have $f(x) = (x - \alpha)f_1(x)$ over F_1 where $f_1(x)$ has degree $n - 1$ and $f'(x) = (x - \beta)f'_1(x)$. The field E is a splitting field for $f_1(x)$ over F_1 : all the roots of $f_1(x)$ are in E and if they were contained in any smaller extension L containing F_1 , then, since F_1 contains α , L would also contain all the roots of $f(x)$, which would contradict the minimality of E as the splitting field of $f(x)$ over F . Similarly E' is a splitting field for $f'_1(x)$ over F'_1 . Since the degrees of $f_1(x)$ and $f'_1(x)$ are less than n , by induction there exists a map $\sigma : E \xrightarrow{\sim} E'$ extending the isomorphism $\sigma' : F_1 \xrightarrow{\sim} F'_1$. This gives the extended diagram:

$$\begin{array}{ccc} \sigma : & E & \xrightarrow{\sim} E' \\ & | & | \\ \sigma' : & F_1 & \xrightarrow{\sim} F'_1 \\ & | & | \\ \varphi : & F & \xrightarrow{\sim} F'. \end{array}$$

Then as the diagram indicates, σ restricted to F_1 is the isomorphism σ' , so in particular σ restricted to F is σ' restricted to F , which is φ , showing that σ is an extension of φ , completing the proof.

Corollary 28. (Uniqueness of Splitting Fields) Any two splitting fields for a polynomial $f(x) \in F[x]$ over a field F are isomorphic.

Proof: Take φ to be the identity mapping from F to itself and E and E' to be two splitting fields for $f(x) (= f'(x))$.

As we mentioned before, this result justifies the terminology of *the* splitting field for $f(x)$ over F , since any two are isomorphic. Splitting fields play a natural role in

the study of algebraic elements (if you are adjoining one root of a polynomial, why not adjoin *all* the roots?) and so take a particularly important role in Galois Theory.

We end this section with a discussion of field extensions of F which contain all the roots of *all* polynomials over F .

Definition. The field \overline{F} is called an *algebraic closure* of F if \overline{F} is algebraic over F and if every polynomial $f(x) \in F[x]$ splits completely over \overline{F} (so that \overline{F} can be said to contain all the elements algebraic over F).

Definition. A field K is said to be *algebraically closed* if every polynomial with coefficients in K has a root in K .

It is not obvious that algebraically closed fields exist nor that there exists an algebraic closure of a given field F (we shall prove this shortly).

Note that if K is algebraically closed, then in fact every $f(x) \in K[x]$ has *all* its roots in K , since by definition $f(x)$ has a root $\alpha \in K$, hence has a factor $x - \alpha$ in $K[x]$. The remaining factor of $f(x)$ then is a polynomial in $K[x]$, hence has a root, so has a linear factor etc., so that $f(x)$ must split completely. Hence if K is algebraically closed, then K itself is an algebraic closure of K and the converse is obvious, so that $K = \overline{K}$ if and only if K is algebraically closed.

The next result shows that the process of “taking the algebraic closure” actually stops after one step — taking the algebraic closure of an algebraic closure does not give a larger field: the field is already algebraically closed (notationally: $\overline{\overline{F}} = \overline{F}$).

Proposition 29. Let \overline{F} be an algebraic closure of F . Then \overline{F} is algebraically closed.

Proof: Let $f(x)$ be a polynomial in $\overline{F}[x]$ and let α be a root of $f(x)$. Then α generates an algebraic extension $\overline{F}(\alpha)$ of \overline{F} , and \overline{F} is algebraic over F . By Theorem 20, $\overline{F}(\alpha)$ is algebraic over F so in particular its element α is algebraic over F . But then $\alpha \in \overline{F}$, showing \overline{F} is algebraically closed.

Given a field F we have already shown how to construct (finite) extensions of F containing all the roots of any given polynomial $f(x) \in F[x]$. Intuitively, an algebraic closure of F is given by the field “generated” by all of these fields. The difficulty with this is “generated” *where?*, since they are not all subfields of a given field. For a *finite* collection of polynomials $f_1(x), \dots, f_k(x)$, we can identify their splitting fields as subfields of the splitting field of the product polynomial $f_1(x) \cdots f_k(x)$, but the same idea used for an *infinite* number of polynomials requires numerous “bookkeeping” identifications and an application of Zorn’s Lemma.

We shall instead construct an algebraic closure of F by first constructing an algebraically closed field containing F . The proof uses a clever idea of Artin which very neatly solves the “bookkeeping” problem of constructing a field containing the appropriate roots of polynomials (which also ultimately relies on Zorn’s Lemma) by introducing a separate variable for every polynomial.

Proposition 30. For any field F there exists an algebraically closed field K containing F .

Proof: For every nonconstant monic polynomial $f = f(x)$ with coefficients in F , let x_f denote an indeterminate and consider the polynomial ring $F[\dots, x_f, \dots]$ generated over F by the variables x_f . In this polynomial ring consider the ideal I generated by the polynomials $f(x_f)$. If this ideal is not proper, then 1 is an element of the ideal, hence we have a relation

$$g_1 f_1(x_{f_1}) + g_2 f_2(x_{f_2}) + \cdots + g_n f_n(x_{f_n}) = 1$$

where the g_i , $i = 1, 2, \dots, n$, are polynomials in the x_f . For $i = 1, 2, \dots, n$ let $x_{f_i} = x_i$ and let x_{n+1}, \dots, x_m be the remaining variables occurring in the polynomials g_j , $j = 1, 2, \dots, n$. Then the relation above reads

$$g_1(x_1, x_2, \dots, x_m) f_1(x_1) + \cdots + g_n(x_1, x_2, \dots, x_m) f_n(x_n) = 1.$$

Let F' be a finite extension of F containing a root α_i of $f_i(x)$ for $i = 1, 2, \dots, n$. Letting $x_i = \alpha_i$, $i = 1, 2, \dots, n$ and setting $x_{n+1} = \cdots = x_m = 0$, say, in the polynomial equation above would imply that $0 = 1$ in F' , clearly impossible.

Since the ideal I is a proper ideal, it is contained in a maximal ideal \mathcal{M} (this is where Zorn's Lemma is used). Then the quotient

$$K_1 = F[\dots, x_f, \dots]/\mathcal{M}$$

is a field containing (an isomorphic copy of) F . Each of the polynomials f has a root in K_1 by construction, namely the image of x_f , since $f(x_f) \in I \subseteq \mathcal{M}$. We have constructed a field K_1 in which every polynomial with coefficients from F has a root. Performing the same construction with K_1 instead of F gives a field K_2 containing K_1 in which all polynomials with coefficients from K_1 have a root. Continuing in this fashion we obtain a sequence of fields

$$F = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_j \subseteq K_{j+1} \subseteq \cdots$$

where every polynomial with coefficients in K_j has a root in K_{j+1} , $j = 0, 1, \dots$. Let

$$K = \bigcup_{j \geq 0} K_j$$

be the union of these fields. Then K is clearly a field containing F . Since K is the union of the fields K_j , the coefficients of any polynomial $h(x)$ in $K[x]$ all lie in some field K_N for N sufficiently large. But then $h(x)$ has a root in K_{N+1} , so has a root in K . It follows that K is algebraically closed, completing the proof.

We now use the algebraically closed field containing F to construct an algebraic closure of F :

Proposition 31. Let K be an algebraically closed field and let F be a subfield of K . Then the collection of elements \overline{F} of K that are algebraic over F is an algebraic closure of F . An algebraic closure of F is unique up to isomorphism.

Proof: By definition, \overline{F} is an algebraic extension of F . Every polynomial $f(x) \in F[x]$ splits completely over K into linear factors $x - \alpha$ (the same is true for every

polynomial even in $K[x]$). But each α is a root of $f(x)$, so is algebraic over F , hence is an element of \overline{F} . It follows that all the linear factors $x - \alpha$ have coefficients in \overline{F} , i.e., $f(x)$ splits completely in $\overline{F}[x]$ and \overline{F} is an algebraic closure of F .

The uniqueness (up to isomorphism) of the algebraic closure is natural in light of the uniqueness (up to isomorphism) of splitting fields, and is proved along the same lines together with an application of Zorn's Lemma and will be omitted.

We shall prove later using Galois theory the following result (purely analytic proofs using complex analysis also exist).

Theorem. (Fundamental Theorem of Algebra) The field \mathbb{C} is algebraically closed.

By Proposition 31, we immediately obtain:

Corollary 32. The field \mathbb{C} contains an algebraic closure for any of its subfields. In particular, $\overline{\mathbb{Q}}$, the collection of complex numbers algebraic over \mathbb{Q} , is an algebraic closure of \mathbb{Q} .

The point of these considerations is that all the computations involving elements algebraic over a field F may be viewed as taking place in one (large) field, namely \overline{F} . Similarly, we can speak sensibly of the composite of any collection of algebraic extensions by viewing them all as subfields of an algebraic closure. In the case of \mathbb{Q} or finite extensions of \mathbb{Q} we may consider all of our computations as occurring in \mathbb{C} .

EXERCISES

1. Determine the splitting field and its degree over \mathbb{Q} for $x^4 - 2$.
2. Determine the splitting field and its degree over \mathbb{Q} for $x^4 + 2$.
3. Determine the splitting field and its degree over \mathbb{Q} for $x^4 + x^2 + 1$.
4. Determine the splitting field and its degree over \mathbb{Q} for $x^6 - 4$.
5. Let K be a finite extension of F . Prove that K is a splitting field over F if and only if every irreducible polynomial in $F[x]$ that has a root in K splits completely in $K[x]$. [Use Theorems 8 and 27.]
6. Let K_1 and K_2 be finite extensions of F contained in the field K , and assume both are splitting fields over F .
 - (a) Prove that their composite K_1K_2 is a splitting field over F .
 - (b) Prove that $K_1 \cap K_2$ is a splitting field over F . [Use the preceding exercise.]

13.5 SEPARABLE AND INSEPARABLE EXTENSIONS

Let F be a field and let $f(x) \in F[x]$ be a polynomial. Over a splitting field for $f(x)$ we have the factorization

$$f(x) = (x - \alpha_1)^{n_1}(x - \alpha_2)^{n_2} \cdots (x - \alpha_k)^{n_k}$$

where $\alpha_1, \alpha_2, \dots, \alpha_k$ are distinct elements of the splitting field and $n_i \geq 1$ for all i . Recall that α_i is called a *multiple* root if $n_i > 1$ and is called a *simple* root if $n_i = 1$. The integer n_i is called the *multiplicity* of the root α_i .

Definition. A polynomial over F is called *separable* if it has no multiple roots (i.e., all its roots are distinct). A polynomial which is not separable is called *inseparable*.

Note that if a polynomial $f(x)$ has distinct roots in one splitting field then $f(x)$ has distinct roots in any splitting field (since this is equivalent to $f(x)$ factoring into distinct linear factors, and there is an isomorphism over F between any two splitting fields of $f(x)$ that is bijective on its roots), so that we need not specify the field containing all the roots of $f(x)$.

Examples

- (1) The polynomial $x^2 - 2$ is separable over \mathbb{Q} since its two roots $\pm\sqrt{2}$ are distinct. The polynomial $(x^2 - 2)^n$ for any $n \geq 2$ is inseparable since it has the multiple roots $\pm\sqrt{2}$, each with multiplicity n .
- (2) The polynomial $x^2 - t$ ($= x^2 + t$) over the field $F = \mathbb{F}_2(t)$ of rational functions in t with coefficients from \mathbb{F}_2 is irreducible as we've seen before, but is not separable. If \sqrt{t} denotes a root in some extension field (note that $\sqrt{t} \notin F$), then

$$(x - \sqrt{t})^2 = x^2 - 2x\sqrt{t} + t = x^2 + t = x^2 - t$$

since F is a field of characteristic 2. Hence this irreducible polynomial has only one root (with multiplicity 2), so is not separable over F .

There is a simple criterion to check whether a polynomial has multiple roots.

Definition. The *derivative* of the polynomial

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0 \in F[x]$$

is defined to be the polynomial

$$D_x f(x) = n a_n x^{n-1} + (n-1) a_{n-1} x^{n-2} + \cdots + 2 a_2 x + a_1 \in F[x].$$

This formula is nothing but the usual formula for the derivative of a polynomial familiar from calculus. It is purely algebraic and so can be applied to a polynomial over an arbitrary field F , where the analytic notion of derivative (involving limits — a *continuous* operation) may not exist.

The usual (calculus) formulas for derivatives hold for derivatives in this situation as well, for example the formulas for the derivative of a sum and of a product:

$$\begin{aligned} D_x(f(x) + g(x)) &= D_x f(x) + D_x g(x) \\ D_x(f(x)g(x)) &= f(x)D_x g(x) + (D_x f(x))g(x). \end{aligned}$$

These formulas can be proved directly from the definition for polynomials and do not require any limiting operations and are left as an exercise.

The next proposition shows that the separability of $f(x)$ can be determined by the Euclidean Algorithm in the field where the coefficients of $f(x)$ lie, without passing to a splitting field and factoring $f(x)$.

Proposition 33. A polynomial $f(x)$ has a multiple root α if and only if α is also a root of $D_x f(x)$, i.e., $f(x)$ and $D_x f(x)$ are both divisible by the minimal polynomial for α . In particular, $f(x)$ is separable if and only if it is relatively prime to its derivative: $(f(x), D_x f(x)) = 1$.

Proof: Suppose first that α is a multiple root of $f(x)$. Then over a splitting field,

$$f(x) = (x - \alpha)^n g(x)$$

for some integer $n \geq 2$ and some polynomial $g(x)$. Taking derivatives we obtain

$$D_x f(x) = n(x - \alpha)^{n-1} g(x) + (x - \alpha)^n D_x g(x)$$

which shows ($n \geq 2$) that $D_x f(x)$ has α as a root.

Conversely, suppose that α is a root of both $f(x)$ and $D_x f(x)$. Then write

$$f(x) = (x - \alpha)h(x)$$

for some polynomial $h(x)$ and take the derivative:

$$D_x f(x) = h(x) + (x - \alpha)D_x h(x).$$

Since $D_x f(\alpha) = 0$ by assumption, substituting α into the last equation shows that $h(\alpha) = 0$. Hence $h(x) = (x - \alpha)h_1(x)$ for some polynomial $h_1(x)$, and

$$f(x) = (x - \alpha)^2 h_1(x)$$

showing that α is a multiple root of $f(x)$.

The equivalence with divisibility by the minimal polynomial for α follows from Proposition 9. The last statement is then clear (let α denote any root of a common factor of $f(x)$ and $D_x f(x)$). \(\checkmark\)

Examples

- (1) The polynomial $x^{p^n} - x$ over \mathbb{F}_p has derivative $p^n x^{p^n-1} - 1 = -1$ since the field has characteristic p . Since in this case the derivative has no roots at all, it follows that the polynomial has no multiple roots, hence is separable.
- (2) The polynomial $x^n - 1$ has derivative nx^{n-1} . Over any field of characteristic not dividing n (including characteristic 0) this polynomial has only the root 0 (of multiplicity $n-1$), which is not a root of $x^n - 1$. Hence $x^n - 1$ is separable and there are n distinct n^{th} roots of unity. We saw this directly over \mathbb{Q} by exhibiting n distinct solutions over \mathbb{C} .
- (3) If F is of characteristic p and p divides n , then there are fewer than n distinct n^{th} roots of unity over F : in this case the derivative is identically 0 since $n = 0$ in F . In fact every root of $x^n - 1$ is multiple in this case.

Corollary 34. Every irreducible polynomial over a field of characteristic 0 (for example, \mathbb{Q}) is separable. A polynomial over such a field is separable if and only if it is the product of distinct irreducible polynomials.

Proof: Suppose F is a field of characteristic 0 and $p(x) \in F[x]$ is irreducible of degree n . Then the derivative $D_x p(x)$ is a polynomial of degree $n-1$. Up to constant factors the only factors of $p(x)$ in $F[x]$ are 1 and $p(x)$, so $D_x p(x)$ must be

relatively prime to $p(x)$. This shows that any irreducible polynomial over a field of characteristic 0 is separable. The second statement of the corollary is then clear since distinct irreducibles never have zeros in common (by Proposition 9).

The point in the proof of the corollary that can fail in characteristic p is the statement that the derivative $D_x p(x)$ is of degree $n - 1$. In characteristic p the derivative of any power x^{pm} of x^p is identically 0:

$$D_x(x^{pm}) = pmx^{pm-1} = 0$$

so it is possible for the degree of the derivative to decrease by more than one. If the derivative $D_x p(x)$ of the *irreducible* polynomial $p(x)$ is nonzero, however, then just as before we conclude that $p(x)$ must be separable.

It is clear from the definition of the derivative that if $p(x)$ is a polynomial whose derivative is 0, then every exponent of x in $p(x)$ must be a multiple of p where p is the characteristic of F :

$$p(x) = a_m x^{mp} + a_{m-1} x^{(m-1)p} + \cdots + a_1 x^p + a_0.$$

Letting

$$p_1(x) = a_m x^m + a_{m-1} x^{m-1} + \cdots + a_1 x + a_0$$

we see that $p(x)$ is a polynomial in x^p , namely $p(x) = p_1(x^p)$.

We now prove a simple but important result about raising to the p^{th} power in a field of characteristic p .

Proposition 35. Let F be a field of characteristic p . Then for any $a, b \in F$,

$$(a + b)^p = a^p + b^p, \quad \text{and} \quad (ab)^p = a^p b^p.$$

Put another way, the p^{th} -power map defined by $\varphi(a) = a^p$ is an injective field homomorphism from F to F .

Proof: The Binomial Theorem for expanding $(a + b)^n$ for any positive integer n holds (by the standard induction proof) over any commutative ring:

$$(a + b)^n = a^n + \binom{n}{1} a^{n-1} b + \cdots + \binom{n}{i} a^{n-i} b^i + \cdots + b^n.$$

It should be observed that the binomial coefficients

$$\binom{n}{i} = \frac{n!}{i!(n-i)!}$$

are integers (recall that $m\alpha$ for $m \in \mathbb{Z}$ is defined for α an element of any ring) and here are elements of the prime field.

If p is a prime, then the binomial coefficients $\binom{p}{i}$ for $i = 1, 2, \dots, p - 1$ are all divisible by p since for these values of i the numbers $i!$ and $(p - i)!$ only involve factors smaller than p , hence are relatively prime to p and so cannot cancel the factor of p in the numerator of the expression $\frac{p!}{i!(p - i)!}$. It follows that over a field of characteristic p all the intermediate terms in the expansion of $(a + b)^p$ are 0, which gives the first equation of the proposition. The second equation is trivial, as is the fact that φ is injective.

Definition. The map in Proposition 35 is called the *Frobenius endomorphism* of F .

Corollary 36. Suppose that \mathbb{F} is a finite field of characteristic p . Then every element of \mathbb{F} is a p^{th} power in \mathbb{F} (notationally, $\mathbb{F} = \mathbb{F}^p$).

Proof: The injectivity of the Frobenius endomorphism of \mathbb{F} implies that it is also surjective when \mathbb{F} is finite, which is the statement of the corollary.

We now prove the analogue of Corollary 34 for finite fields.

Let \mathbb{F} be a finite field and suppose that $p(x) \in \mathbb{F}[x]$ is an irreducible polynomial with coefficients in \mathbb{F} . If $p(x)$ were inseparable then we have seen that $p(x) = q(x^p)$ for some polynomial $q(x) \in \mathbb{F}[x]$. Let

$$q(x) = a_m x^m + a_{m-1} x^{m-1} + \cdots + a_1 x + a_0.$$

By Corollary 36, each a_i , $i = 1, 2, \dots, m$ is a p^{th} power in \mathbb{F} , say $a_i = b_i^p$. Then by Proposition 35 we have

$$\begin{aligned} p(x) &= q(x^p) = a_m (x^p)^m + a_{m-1} (x^p)^{m-1} + \cdots + a_1 x^p + a_0 \\ &= b_m^p (x^p)^m + b_{m-1}^p (x^p)^{m-1} + \cdots + b_1^p x^p + b_0^p \\ &= (b_m x^m)^p + (b_{m-1} x^{m-1})^p + \cdots + (b_1 x)^p + (b_0)^p \\ &= (b_m x^m + b_{m-1} x^{m-1} + \cdots + b_1 x + b_0)^p \end{aligned}$$

which shows that $p(x)$ is the p^{th} power of a polynomial in $\mathbb{F}[x]$, a contradiction to the irreducibility of $p(x)$. This proves:

Proposition 37. Every irreducible polynomial over a finite field \mathbb{F} is separable. A polynomial in $\mathbb{F}[x]$ is separable if and only if it is the product of distinct irreducible polynomials in $\mathbb{F}[x]$.

The important part of the proof of this result is the fact that every element in the characteristic p field \mathbb{F} was a p^{th} power in \mathbb{F} . This suggests the following definition:

Definition. A field K of characteristic p is called *perfect* if every element of K is a p^{th} power in K , i.e., $K = K^p$. Any field of characteristic 0 is also called perfect.

With this definition, we see that we have proved that every irreducible polynomial over a perfect field is separable. It is not hard to see that if K is not perfect then there are inseparable irreducible polynomials.

Example: (Existence and Uniqueness of Finite Fields)

Let $n > 0$ be any positive integer and consider the splitting field of the polynomial $x^{p^n} - x$ over \mathbb{F}_p . We have already seen that this polynomial is separable, hence has precisely p^n roots. Let α and β be any two roots of this polynomial, so that $\alpha^{p^n} = \alpha$ and $\beta^{p^n} = \beta$. Then $(\alpha\beta)^{p^n} = \alpha\beta$, $(\alpha^{-1})^{p^n} = \alpha^{-1}$ and by Proposition 35 also

$$(\alpha + \beta)^{p^n} = \alpha^{p^n} + \beta^{p^n} = \alpha + \beta.$$

Hence the set \mathbb{F} consisting of the p^n distinct roots of $x^{p^n} - x$ over \mathbb{F}_p is *closed* under addition, multiplication and inverses in its splitting field. It follows that \mathbb{F} is a subfield, hence in fact must be the splitting field. Since the number of elements is p^n , we have $[\mathbb{F} : \mathbb{F}_p] = n$, which shows that there exist finite fields of degree n over \mathbb{F}_p for any $n > 0$.

Let now \mathbb{F} be any finite field of characteristic p . If \mathbb{F} is of dimension n over its prime subfield \mathbb{F}_p , then \mathbb{F} has precisely p^n elements. Since the multiplicative group \mathbb{F}^\times is (in fact cyclic) of order $p^n - 1$, we have $\alpha^{p^n-1} = 1$ for every $\alpha \neq 0$ in \mathbb{F} , so that $\alpha^{p^n} = \alpha$ for every $\alpha \in \mathbb{F}$. But this means α is a root of $x^{p^n} - x$, hence \mathbb{F} is contained in a splitting field for this polynomial. Since we have seen that the splitting field has order p^n this shows that \mathbb{F} is a splitting field for $x^{p^n} - x$. Since splitting fields are unique up to isomorphism, this proves that *finite fields of any order p^n exist and are unique up to isomorphism*. We shall denote the finite field of order p^n by \mathbb{F}_{p^n} .

We shall consider finite fields more later.

We now investigate further the structure of inseparable irreducible polynomials over fields of characteristic p . We have seen above that if $p(x)$ is an irreducible polynomial which is not separable, then its derivative $D_x p(x)$ is identically 0, so that $p(x) = p_1(x^p)$ for some polynomial $p_1(x)$. The polynomial $p_1(x)$ may or may not itself be separable. If not, then it too is a polynomial in x^p , $p_1(x) = p_2(x^p)$, so that $p(x)$ is a polynomial in x^{p^2} : $p(x) = p_2(x^{p^2})$. Continuing in this fashion we see that there is a uniquely defined power p^k of p such that $p(x) = p_k(x^{p^k})$ where $p_k(x)$ has nonzero derivative. It is clear that $p_k(x)$ is irreducible since any factorization of $p_k(x)$ would, after replacing x by x^{p^k} , immediately imply a factorization of the irreducible $p(x)$. It follows that $p_k(x)$ is separable. We summarize this as:

Proposition 38. Let $p(x)$ be an irreducible polynomial over a field F of characteristic p . Then there is a unique integer $k \geq 0$ and a unique irreducible separable polynomial $p_{sep}(x) \in F[x]$ such that

$$p(x) = p_{sep}(x^{p^k}).$$

Definition. Let $p(x)$ be an irreducible polynomial over a field of characteristic p . The degree of $p_{sep}(x)$ in the last proposition is called the *separable degree* of $p(x)$, denoted $\deg_s p(x)$. The integer p^k in the proposition is called the *inseparable degree* of $p(x)$, denoted $\deg_i p(x)$.

From the definitions and the proposition we see that $p(x)$ is separable if and only if its inseparability degree is 1 if and only if its degree is equal to its separable degree. Also, computing degrees in the relation $p(x) = p_{sep}(x^{p^k})$ we see that

$$\deg p(x) = \deg_s p(x) \deg_i p(x).$$

Examples

- (1) The polynomial $p(x) = x^2 - t$ over $F = \mathbb{F}_2(t)$ considered above has derivative 0, hence is not separable (as we determined earlier). Here $p_{sep}(x) = x - t$ with inseparability degree 2.

- (2) The polynomial $p(x) = x^{2^m} - t$ over $F = \mathbb{F}_2(t)$ is irreducible with the same separable polynomial part, but with inseparability degree 2^m .
- (3) The polynomial $(x^{p^2} - t)(x^p - t)$ over $F = \mathbb{F}_p(t)$ has (two) inseparable irreducible factors so is inseparable. This polynomial cannot be written in the form $f_{sep}(x^{p^k})$ where $f_{sep}(x)$ is separable, which is the reason we restricted to *irreducible* polynomials above. This example also shows that there is no analogous factorization to define the separable and inseparable degrees of a general polynomial.

The notion of separability carries over to the fields generated by the roots of these polynomials.

Definition. The field K is said to be *separable* (or *separably algebraic*) over F if every element of K is the root of a separable polynomial over F (equivalently, the minimal polynomial over F of every element of K is separable). A field which is not separable is *inseparable*.

We have seen that the issue of separability is straightforward for finite extensions of perfect fields since for these fields the minimal polynomial of an algebraic element is irreducible hence separable.

Corollary 39. Every finite extension of a perfect field is separable. In particular, every finite extension of either \mathbb{Q} or a finite field is separable.

We shall consider separable and inseparable extensions more after developing some Galois Theory, in particular defining the separable and inseparable *degree* of the extension K/F .

EXERCISES

1. Prove that the derivative D_x of a polynomial satisfies $D_x(f(x) + g(x)) = D_x(f(x)) + D_x(g(x))$ and $D_x(f(x)g(x)) = D_x(f(x))g(x) + D_x(g(x))f(x)$ for any two polynomials $f(x)$ and $g(x)$.
2. Find all irreducible polynomials of degrees 1, 2 and 4 over \mathbb{F}_2 and prove that their product is $x^{16} - x$.
3. Prove that d divides n if and only if $x^d - 1$ divides $x^n - 1$. [Note that if $n = qd + r$ then $x^n - 1 = (x^{qd+r} - x^r) + (x^r - 1)$.]
4. Let $a > 1$ be an integer. Prove for any positive integers n, d that d divides n if and only if $a^d - 1$ divides $a^n - 1$ (cf. the previous exercise). Conclude in particular that $\mathbb{F}_{p^d} \subseteq \mathbb{F}_{p^n}$ if and only if d divides n .
5. For any prime p and any nonzero $a \in \mathbb{F}_p$ prove that $x^p - x + a$ is irreducible and separable over \mathbb{F}_p . [For the irreducibility: One approach — prove first that if α is a root then $\alpha + 1$ is also a root. Another approach — suppose it's reducible and compute derivatives.]
6. Prove that $x^{p^n-1} - 1 = \prod_{\alpha \in \mathbb{F}_{p^n}^\times} (x - \alpha)$. Conclude that $\prod_{\alpha \in \mathbb{F}_{p^n}^\times} \alpha = (-1)^{p^n}$ so the product of the nonzero elements of a finite field is $+1$ if $p = 2$ and -1 if p is odd. For p odd and $n = 1$ derive *Wilson's Theorem*: $(p - 1)! \equiv -1 \pmod{p}$.

7. Suppose K is a field of characteristic p which is not a perfect field: $K \neq K^p$. Prove there exist irreducible inseparable polynomials over K . Conclude that there exist inseparable finite extensions of K .

8. Prove that $f(x)^p = f(x^p)$ for any polynomial $f(x) \in \mathbb{F}_p[x]$.

9. Show that the binomial coefficient $\binom{pn}{pi}$ is the coefficient of x^{pi} in the expansion of $(1+x)^{pn}$.

Working over \mathbb{F}_p show that this is the coefficient of $(x^p)^i$ in $(1+x^p)^n$ and hence prove that $\binom{pn}{pi} \equiv \binom{n}{i} \pmod{p}$.

10. Let $f(x_1, x_2, \dots, x_n) \in \mathbb{Z}[x_1, x_2, \dots, x_n]$ be a polynomial in the variables x_1, x_2, \dots, x_n with integer coefficients. For any prime p prove that the polynomial

$$f(x_1, x_2, \dots, x_n)^p - f(x_1^p, x_2^p, \dots, x_n^p) \in \mathbb{Z}[x_1, x_2, \dots, x_n]$$

has all its coefficients divisible by p .

11. Suppose $K[x]$ is a polynomial ring over the field K and F is a subfield of K . If F is a perfect field and $f(x) \in F[x]$ has no repeated irreducible factors in $F[x]$, prove that $f(x)$ has no repeated irreducible factors in $K[x]$.

13.6 CYCLOTOMIC POLYNOMIALS AND EXTENSIONS

The purpose of this section is to prove that the cyclotomic extension

$$\mathbb{Q}(\zeta_n)/\mathbb{Q}$$

generated by the n^{th} roots of unity over \mathbb{Q} introduced in Section 4 is of degree $\varphi(n)$ where φ denotes Euler's phi-function (= the number of integers a , $1 \leq a < n$ relatively prime to n = the order of the group $(\mathbb{Z}/n\mathbb{Z})^\times$).

Definition. Let μ_n denote the group of n^{th} roots of unity over \mathbb{Q} .

Then as we have already observed, $\mathbb{Z}/n\mathbb{Z} \cong \mu_n$ as groups (under multiplication on the right, addition on the left), given explicitly by the map $a \mapsto (\zeta_n)^a$ for a fixed primitive n^{th} root of unity. The primitive n^{th} roots of unity are given by the residue classes prime to n so there are precisely $\varphi(n)$ primitive n^{th} roots of unity.

If d is a divisor of n and ζ is a d^{th} root of unity, then ζ is also an n^{th} root of unity since $\zeta^n = (\zeta^d)^{n/d} = 1$. Hence

$$\mu_d \subseteq \mu_n \quad \text{for all } d \mid n.$$

Conversely, the order of any element of the group μ_n is a divisor of n so that if ζ is an n^{th} root of unity which is also a d^{th} root of unity for some smaller d then $d \mid n$.

Definition. Define the n^{th} cyclotomic polynomial $\Phi_n(x)$ to be the polynomial whose roots are the primitive n^{th} roots of unity:

$$\Phi_n(x) = \prod_{\zeta \text{ primitive } \in \mu_n} (x - \zeta) = \prod_{\substack{1 \leq a < n \\ (a, n) = 1}} (x - \zeta_n^a)$$

(which is of degree $\varphi(n)$).

The roots of the polynomial $x^n - 1$ are precisely the n^{th} roots of unity so we have the factorization

$$x^n - 1 = \prod_{\substack{\zeta^n = 1 \\ \text{i.e. } \zeta \in \mu_n}} (x - \zeta).$$

If we group together the factors $(x - \zeta)$ where ζ is an element of order d in μ_n (i.e., ζ is a primitive d^{th} root of unity) we obtain

$$x^n - 1 = \prod_{d|n} \prod_{\substack{\zeta \in \mu_d \\ \zeta \text{ primitive}}} (x - \zeta).$$

The inner product is $\Phi_d(x)$ by definition so we have the factorization

$$x^n - 1 = \prod_{d|n} \Phi_d(x). \quad (13.4)$$

Note incidentally that comparing degrees gives the identity

$$n = \sum_{d|n} \varphi(d).$$

This factorization allows us to compute $\Phi_n(x)$ for any n recursively: clearly $\Phi_1(x) = x - 1$ and $\Phi_2(x) = x + 1$. Then

$$x^3 - 1 = \Phi_1(x)\Phi_3(x) = (x - 1)\Phi_3(x)$$

which gives

$$\Phi_3(x) = x^2 + x + 1.$$

Similarly

$$x^4 - 1 = \Phi_1(x)\Phi_2(x)\Phi_4(x) = (x - 1)(x + 1)\Phi_4(x)$$

gives

$$\Phi_4(x) = x^2 + 1$$

(in these cases these could also be obtained directly from the explicit roots of unity). Continuing in this fashion we can compute $\Phi_n(x)$ for any n . Note also that for p a prime we recover our polynomial

$$\Phi_p(x) = x^{p-1} + x^{p-2} + \cdots + x + 1.$$

For some small values of n the polynomials are

$$\Phi_5(x) = x^4 + x^3 + x^2 + x + 1$$

$$\Phi_6(x) = x^2 - x + 1$$

$$\Phi_7(x) = x^6 + x^5 + x^4 + x^3 + x^2 + x + 1$$

$$\Phi_8(x) = x^4 + 1$$

$$\Phi_9(x) = x^6 + x^3 + 1$$

$$\Phi_{10}(x) = x^4 - x^3 + x^2 - x + 1$$

$$\Phi_{11}(x) = x^{10} + x^9 + \cdots + x + 1$$

$$\Phi_{12}(x) = x^4 - x^2 + 1.$$

For all the values computed above, $\Phi_n(x)$ was a (monic) polynomial with integer coefficients. This is always the case:

Lemma 40. The cyclotomic polynomial $\Phi_n(x)$ is a monic polynomial in $\mathbb{Z}[x]$ of degree $\varphi(n)$.

Proof: It is clear that $\Phi_n(x)$ is monic and has degree $\varphi(n)$. We must show the coefficients lie in \mathbb{Z} . We use induction on n . The result is true for $n = 1$ (and $n \leq 12$). Assume by induction that $\Phi_d(x) \in \mathbb{Z}[x]$ for all $1 \leq d < n$. Then $x^n - 1 = f(x)\Phi_n(x)$ where $f(x) = \prod_{d|n} \Phi_d(x)$ is monic and has coefficients in \mathbb{Z} . Since $f(x)$ clearly divides $x^n - 1$ in $F[x]$ where $F = \mathbb{Q}(\zeta_n)$ is the field of n^{th} roots of unity and both $f(x)$ and $x^n - 1$ have coefficients in \mathbb{Q} , $f(x)$ divides $x^n - 1$ in $\mathbb{Q}[x]$ by the Division Algorithm (cf. the remark at the end of Section 9.2). By Gauss' Lemma, $f(x)$ divides $x^n - 1$ in $\mathbb{Z}[x]$, hence $\Phi_n(x) \in \mathbb{Z}[x]$.

We remark in passing that while all the coefficients of $\Phi_n(x)$ in the examples computed above were 0, ± 1 , it is known that there are cyclotomic polynomials with arbitrarily large coefficients.

Theorem 41. The cyclotomic polynomial $\Phi_n(x)$ is an irreducible monic polynomial in $\mathbb{Z}[x]$ of degree $\varphi(n)$.

Proof: We must show that $\Phi_n(x)$ is irreducible. If not then we have a factorization

$$\Phi_n(x) = f(x)g(x) \quad \text{with } f(x), g(x) \text{ monic in } \mathbb{Z}[x]$$

where we take $f(x)$ to be an *irreducible* factor of $\Phi_n(x)$. Let ζ be a primitive n^{th} root of 1 which is a root of $f(x)$ (so then $f(x)$ is the minimal polynomial for ζ over \mathbb{Q}) and let p denote *any* prime not dividing n . Then ζ^p is again a primitive n^{th} root of 1, hence is a root of either $f(x)$ or $g(x)$.

Suppose $g(\zeta^p) = 0$. Then ζ is a root of $g(x^p)$ and since $f(x)$ is the minimal polynomial for ζ , $f(x)$ must divide $g(x^p)$ in $\mathbb{Z}[x]$, say

$$g(x^p) = f(x)h(x), \quad h(x) \in \mathbb{Z}[x].$$

If we reduce this equation mod p , we obtain

$$\bar{g}(x^p) = \bar{f}(x)\bar{h}(x) \quad \text{in } \mathbb{F}_p[x].$$

By the remarks of the last section,

$$\bar{g}(x^p) = (\bar{g}(x))^p$$

so we have the equation

$$(\bar{g}(x))^p = \bar{f}(x)\bar{h}(x)$$

in the U.F.D. $\mathbb{F}_p[x]$. It follows that $\bar{f}(x)$ and $\bar{g}(x)$ have a factor in common in $\mathbb{F}_p[x]$.

Now, from $\Phi_n(x) = f(x)g(x)$ we see by reducing mod p that $\bar{\Phi}_n(x) = \bar{f}(x)\bar{g}(x)$, and so by the above it follows that $\bar{\Phi}_n(x) \in \mathbb{F}_p[x]$ has a multiple root. But then also $x^n - 1$ would have a multiple root over \mathbb{F}_p since it has $\bar{\Phi}_n(x)$ as a factor. This is a

contradiction since we have seen in the last section that there are n distinct roots of $x^n - 1$ over any field of characteristic not dividing n .

Hence ζ^p must be a root of $f(x)$. Since this applies to every root ζ of $f(x)$, it follows that ζ^a is a root of $f(x)$ for every integer a relatively prime to n : write $a = p_1 p_2 \cdots p_k$ as a product of (not necessarily distinct) primes not dividing n so that ζ^{p_1} is a root of $f(x)$, so also $(\zeta^{p_1})^{p_2}$ is a root of $f(x)$, etc. But this means that every primitive n^{th} root of unity is a root of $f(x)$, i.e., $f(x) = \Phi_n(x)$, showing $\Phi_n(x)$ is irreducible.

Corollary 42. The degree over \mathbb{Q} of the cyclotomic field of n^{th} roots of unity is $\varphi(n)$:

$$[\mathbb{Q}(\zeta_n) : \mathbb{Q}] = \varphi(n).$$

Proof: By the theorem, $\Phi_n(x)$ is the minimal polynomial for any primitive n^{th} root of unity ζ_n .

Example

The cyclotomic field $\mathbb{Q}(\zeta_8)$ of the 8^{th} roots of unity is of degree $\varphi(8) = 4$ over \mathbb{Q} . This field contains the 4^{th} roots of unity, i.e., $\mathbb{Q}(i) \subset \mathbb{Q}(\zeta_8)$ as well as the element $\zeta_8 + \zeta_8^{-1} = \sqrt{2}$ (recall the explicit roots of unity in Section 4). It follows that

$$\mathbb{Q}(\zeta_8) = \mathbb{Q}(i, \sqrt{2}).$$

One interesting number-theoretic application of the cyclotomic polynomials outlined in the exercises is the proof that for any n there are infinitely many primes which are congruent to 1 modulo n . The complete factorization in $\mathbb{F}_p[x]$ of $\Phi_\ell(x)$ for a prime ℓ (which is irreducible in $\mathbb{Z}[x]$) is described in Exercise 8 below.

We shall return to the example of cyclotomic fields after we have developed some Galois Theory.

EXERCISES

- Suppose m and n are relatively prime positive integers. Let ζ_m be a primitive m^{th} root of unity and let ζ_n be a primitive n^{th} root of unity. Prove that $\zeta_m \zeta_n$ is a primitive mn^{th} root of unity.
- Let ζ_n be a primitive n^{th} root of unity and let d be a divisor of n . Prove that ζ_n^d is a primitive $(n/d)^{\text{th}}$ root of unity.
- Prove that if a field contains the n^{th} roots of unity for n odd then it also contains the $2n^{\text{th}}$ roots of unity.
- Prove that if $n = p^k m$ where p is a prime and m is relatively prime to p then there are precisely m distinct n^{th} roots of unity over a field of characteristic p .
- Prove there are only a finite number of roots of unity in any finite extension K of \mathbb{Q} .
- Prove that for n odd, $n > 1$, $\Phi_{2n}(x) = \Phi_n(-x)$.
- Use the Möbius Inversion formula indicated in Section 14.3 to prove

$$\Phi_m(x) = \prod_{d|n} (x^d - 1)^{\mu(m/d)}.$$

8. Let ℓ be a prime and let $\Phi_\ell(x) = \frac{x^\ell - 1}{x - 1} = x^{\ell-1} + x^{\ell-2} + \dots + x + 1 \in \mathbb{Z}[x]$ be the ℓ^{th} cyclotomic polynomial, which is irreducible over \mathbb{Z} by Theorem 41. This exercise determines the factorization of $\Phi_\ell(x)$ modulo p for any prime p . Let ζ denote any fixed primitive ℓ^{th} root of unity.
- (a) Show that if $p = \ell$ then $\Phi_\ell(x) = (x - 1)^{\ell-1} \in \mathbb{F}_\ell[x]$.
 - (b) Suppose $p \neq \ell$ and let f denote the order of p mod ℓ , i.e., f is the smallest power of p with $p^f \equiv 1 \pmod{\ell}$. Use the fact that $\mathbb{F}_{p^n}^\times$ is a cyclic group to show that $n = f$ is the smallest power p^n of p with $\zeta \in \mathbb{F}_{p^n}$. Conclude that the minimal polynomial of ζ over \mathbb{F}_p has degree f .
 - (c) Show that $\mathbb{F}_p(\zeta) = \mathbb{F}_p(\zeta^a)$ for any integer a not divisible by ℓ . [One inclusion is obvious. For the other, note that $\zeta = (\zeta^a)^b$ where b is the multiplicative inverse of a mod ℓ .] Conclude using (b) that, in $\mathbb{F}_p[x]$, $\Phi_\ell(x)$ is the product of $\frac{\ell-1}{f}$ distinct irreducible polynomials of degree f .
 - (d) In particular, prove that, viewed in $\mathbb{F}_p[x]$, $\Phi_7(x) = x^6 + x^5 + \dots + x + 1$ is $(x - 1)^6$ for $p = 7$, a product of distinct linear factors for $p \equiv 1 \pmod{7}$, a product of 3 irreducible quadratics for $p \equiv 6 \pmod{7}$, a product of 2 irreducible cubics for $p \equiv 2, 4 \pmod{7}$, and is irreducible for $p \equiv 3, 5 \pmod{7}$.
9. Suppose A is an $n \times n$ matrix over \mathbb{C} for which $A^k = I$ for some integer $k \geq 1$. Show that A can be diagonalized. Show that the matrix $A = \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix}$ where α is an element of a field of characteristic p satisfies $A^p = I$ and cannot be diagonalized if $\alpha \neq 0$.
10. Let φ denote the Frobenius map $x \mapsto x^p$ on the finite field \mathbb{F}_{p^n} . Prove that φ gives an isomorphism of \mathbb{F}_{p^n} to itself (such an isomorphism is called an *automorphism*). Prove that φ^n is the identity map and that no lower power of φ is the identity.
11. Let φ denote the Frobenius map $x \mapsto x^p$ on the finite field \mathbb{F}_{p^n} as in the previous exercise. Determine the rational canonical form over \mathbb{F}_p for φ considered as an \mathbb{F}_p -linear transformation of the n -dimensional \mathbb{F}_p -vector space \mathbb{F}_{p^n} .
12. Let φ denote the Frobenius map $x \mapsto x^p$ on the finite field \mathbb{F}_{p^n} as in the previous exercise. Determine the Jordan canonical form (over a field containing all the eigenvalues) for φ considered as an \mathbb{F}_p -linear transformation of the n -dimensional \mathbb{F}_p -vector space \mathbb{F}_{p^n} .
13. (*Wedderburn's Theorem on Finite Division Rings*) This exercise outlines a proof (following Witt) of Wedderburn's Theorem that a finite division ring D is a field (i.e., is commutative).
 - (a) Let Z denote the center of D (i.e., the elements of D which commute with every element of D). Prove that Z is a field containing \mathbb{F}_p for some prime p . If $Z = \mathbb{F}_q$ prove that D has order q^n for some integer n [D is a vector space over Z].
 - (b) The nonzero elements D^\times of D form a multiplicative group. For any $x \in D^\times$ show that the elements of D which commute with x form a division ring which contains Z . Show that this division ring is of order q^{m_i} for some integer m_i and that $m_i < n$ if x is not an element of Z .
 - (c) Show that the class equation (Theorem 4.7) for the group D^\times is

$$q^n - 1 = (q - 1) + \sum_{i=1}^r \frac{q^{m_i} - 1}{|C_{D^\times}(x_i)|}$$

where x_1, x_2, \dots, x_r are representatives of the distinct conjugacy classes in D^\times not contained in the center of D^\times . Conclude from (b) that for each i , $|C_{D^\times}(x_i)| = q^{m_i} - 1$ for some $m_i < n$.

- (d) Prove that since $\frac{q^n - 1}{q^{m_i} - 1}$ is an integer (namely, the index $|D^\times : C_{D^\times}(x_i)|$) then m_i divides n (cf. Exercise 4 of Section 5). Conclude that $\Phi_n(x)$ divides $(x^n - 1)/(x^{m_i} - 1)$ and hence that the integer $\Phi_n(q)$ divides $(q^n - 1)/(q^{m_i} - 1)$ for $i = 1, 2, \dots, r$.
- (e) Prove that (c) and (d) imply that $\Phi_n(q) = \prod_{\zeta \text{ primitive}} (q - \zeta)$ divides $q - 1$. Prove that $|q - \zeta| > q - 1$ (complex absolute value) for any root of unity $\zeta \neq 1$ [note that 1 is the closest point on the unit circle in \mathbb{C} to the point q on the real line]. Conclude that $n = 1$, i.e., that $D = \mathbb{Z}$ is a field.

The following exercises provide a proof that for any positive integer m there are infinitely many primes p with $p \equiv 1 \pmod{m}$. This is a special case of *Dirichlet's Theorem on Primes in Arithmetic Progressions* which states more generally that there are infinitely many primes p with $p \equiv a \pmod{m}$ for any a relatively prime to m .

14. Given any monic polynomial $P(x) \in \mathbb{Z}[x]$ of degree at least one show that there are infinitely many distinct prime divisors of the integers

$$P(1), P(2), P(3), \dots, P(n), \dots$$

[Suppose p_1, p_2, \dots, p_k are the only primes dividing the values $P(n)$, $n = 1, 2, \dots$. Let N be an integer with $P(N) = a \neq 0$. Show that $Q(x) = a^{-1} P(N + a p_1 p_2 \dots p_k x)$ is an element of $\mathbb{Z}[x]$ and that $Q(n) \equiv 1 \pmod{p_1 p_2 \dots p_k}$ for $n = 1, 2, \dots$. Conclude that there is some integer M such that $Q(M)$ has a prime factor different from p_1, p_2, \dots, p_k and hence that $P(N + a p_1 p_2 \dots p_k M)$ has a prime factor different from p_1, p_2, \dots, p_k .]

15. Let p be an odd prime not dividing m and let $\Phi_m(x)$ be the m^{th} cyclotomic polynomial. Suppose $a \in \mathbb{Z}$ satisfies $\Phi_m(a) \equiv 0 \pmod{p}$. Prove that a is relatively prime to p and that the order of a in $(\mathbb{Z}/p\mathbb{Z})^\times$ is precisely m . [Since

$$x^m - 1 = \prod_{d|m} \Phi_d(x) = \Phi_m(x) \prod_{\substack{d|m \\ d < m}} \Phi_d(x)$$

we see first that $a^m - 1 \equiv 0 \pmod{p}$ i.e., $a^m \equiv 1 \pmod{p}$. If the order of $a \pmod{p}$ were less than m , then $a^d \equiv 1 \pmod{p}$ for some d dividing m , so then $\Phi_d(a) \equiv 0 \pmod{p}$ for some $d < m$. But then $x^m - 1$ would have a as a multiple root mod p , a contradiction.]

16. Let $a \in \mathbb{Z}$. Show that if p is an odd prime dividing $\Phi_m(a)$ then either p divides m or $p \equiv 1 \pmod{m}$.

17. Prove there are infinitely many primes p with $p \equiv 1 \pmod{m}$.

CHAPTER 14

Galois Theory

14.1 BASIC DEFINITIONS

In the previous chapter we proved the existence of a finite extension of a field F which contains all the roots of a given polynomial $f(x)$ whose coefficients are in F . The main idea of Galois Theory (named for Évariste Galois, 1811–1832) is to consider the relation of the group of permutations of the roots of $f(x)$ to the algebraic structure of its splitting field. The connection is given by the Fundamental Theorem of the next section. It can be viewed as another (extremely elegant) application of the important idea in mathematics that one (in our case algebraic) object *acting* on another provides structural information about both.

In this section we introduce the terminology and basic properties of the objects of interest. Let K be a field.

Definition.

- (1) An isomorphism σ of K with itself is called an *automorphism* of K . The collection of automorphisms of K is denoted $\text{Aut}(K)$. If $\alpha \in K$ we shall write $\sigma\alpha$ for $\sigma(\alpha)$.
- (2) An automorphism $\sigma \in \text{Aut}(K)$ is said to *fix* an element $\alpha \in K$ if $\sigma\alpha = \alpha$. If F is a subset of K (for example, a subfield), then an automorphism σ is said to *fix* F if it fixes all the elements of F , i.e., $\sigma a = a$ for all $a \in F$.

Note that any field has at least one automorphism, the identity map, denoted by 1 and sometimes called the *trivial* automorphism.

The prime field of K is generated by $1 \in K$ and since any automorphism σ takes 1 to 1 (and 0 to 0), i.e., $\sigma(1) = 1$, it follows that $\sigma a = a$ for all a in the prime field. Hence any automorphism of a field K fixes its prime subfield. In particular we see that \mathbb{Q} and \mathbb{F}_p have only the trivial automorphism: $\text{Aut}(\mathbb{Q}) = \{1\}$ and $\text{Aut}(\mathbb{F}_p) = \{1\}$.

Definition. Let K/F be an extension of fields. Let $\text{Aut}(K/F)$ be the collection of automorphisms of K which fix F .

Note that if F is the prime subfield of K then $\text{Aut}(K) = \text{Aut}(K/F)$ since every automorphism of K automatically fixes F .

If σ and τ are automorphisms of K then the composite $\sigma\tau$ (and also the composite $\tau\sigma$, which may not be the same) is defined and is again an automorphism of K .

Proposition 1. $\text{Aut}(K)$ is a group under composition and $\text{Aut}(K/F)$ is a subgroup.

Proof: It is clear that $\text{Aut}(K)$ is a group. If σ and τ are automorphisms of K which fix F then also $\sigma\tau$ and σ^{-1} are the identity on F , which shows that $\text{Aut}(K/F)$ is a subgroup.

The following proposition is extremely useful for determining the automorphisms of algebraic extensions.

Proposition 2. Let K/F be a field extension and let $\alpha \in K$ be algebraic over F . Then for any $\sigma \in \text{Aut}(K/F)$, $\sigma\alpha$ is a root of the minimal polynomial for α over F i.e., $\text{Aut}(K/F)$ permutes the roots of irreducible polynomials. Equivalently, any polynomial with coefficients in F having α as a root also has $\sigma\alpha$ as a root.

Proof: Suppose α satisfies the equation

$$\alpha^n + a_{n-1}\alpha^{n-1} + \cdots + a_1\alpha + a_0 = 0$$

where a_0, a_1, \dots, a_{n-1} are elements of F . Applying the automorphism σ we obtain (using the fact that σ is an additive homomorphism)

$$\sigma(\alpha^n) + \sigma(a_{n-1}\alpha^{n-1}) + \cdots + \sigma(a_1\alpha) + \sigma(a_0) = \sigma(0) = 0.$$

Using the fact that σ is also a multiplicative homomorphism this becomes

$$(\sigma(\alpha))^n + \sigma(a_{n-1})(\sigma(\alpha))^{n-1} + \cdots + \sigma(a_1)(\sigma(\alpha)) + \sigma(a_0) = 0.$$

By assumption, σ fixes all the elements of F , so $\sigma(a_i) = a_i$, $i = 0, 1, \dots, n-1$. Hence

$$(\sigma\alpha)^n + a_{n-1}(\sigma\alpha)^{n-1} + \cdots + a_1(\sigma\alpha) + a_0 = 0.$$

But this says precisely that $\sigma\alpha$ is a root of the same polynomial over F as α . This proves the proposition.

Examples

- (1) Let $K = \mathbb{Q}(\sqrt{2})$. If $\tau \in \text{Aut}(\mathbb{Q}(\sqrt{2})) = \text{Aut}(\mathbb{Q}(\sqrt{2})/\mathbb{Q})$, then $\tau(\sqrt{2}) = \pm\sqrt{2}$ since these are the two roots of the minimal polynomial for $\sqrt{2}$. Since τ fixes \mathbb{Q} , this determines τ completely:

$$\tau(a + b\sqrt{2}) = a \pm b\sqrt{2}.$$

The map $\sqrt{2} \mapsto \sqrt{2}$ is just the identity automorphism 1 of $\mathbb{Q}(\sqrt{2})$. The map $\sigma : \sqrt{2} \mapsto -\sqrt{2}$ is the isomorphism considered in Example 2 following Corollary 13.7. Hence $\text{Aut}(\mathbb{Q}(\sqrt{2})) = \text{Aut}(\mathbb{Q}(\sqrt{2})/\mathbb{Q}) = \{1, \sigma\}$ is a cyclic group of order 2 generated by σ .

- (2) Let $K = \mathbb{Q}(\sqrt[3]{2})$. As before, if $\tau \in \text{Aut}(K/\mathbb{Q})$, then τ is completely determined by its action on $\sqrt[3]{2}$ since

$$\tau(a + b\sqrt[3]{2} + c(\sqrt[3]{2})^2) = a + b\tau\sqrt[3]{2} + c(\tau\sqrt[3]{2})^2.$$

Since $\tau\sqrt[3]{2}$ must be a root of $x^3 - 2$ and the other two roots of this equation are not elements of K (recall the splitting field of this polynomial is degree 6 over \mathbb{Q}), the only possibility is $\tau\sqrt[3]{2} = \sqrt[3]{2}$ i.e., $\tau = 1$. Hence $\text{Aut}(\mathbb{Q}(\sqrt[3]{2})/\mathbb{Q}) = 1$ is the trivial group.

In general, if K is generated over F by some collection of elements, then any automorphism $\sigma \in \text{Aut}(K/F)$ is completely determined by what it does to the generators. If K/F is finite then K is finitely generated over F by algebraic elements so by the proposition the number of automorphisms of K fixing F is finite, i.e., $\text{Aut}(K/F)$ is a finite group. In particular, the automorphisms of a finite extension can be considered as permutations of the roots of a finite number of equations (not every permutation gives rise to an automorphism, however, as Example 2 above illustrates). It was the investigation of permutations of the roots of equations that led Galois to the theory we are describing.

We have associated to each field extension K/F (equivalently, with a subfield F of K) a *group*, $\text{Aut}(K/F)$, the group of automorphisms of K which fix F . One can also reverse this process and associate to each group of automorphisms a field extension.

Proposition 3. Let $H \leq \text{Aut}(K)$ be a subgroup of the group of automorphisms of K . Then the collection F of elements of K fixed by all the elements of H is a subfield of K .

Proof: Let $h \in H$ and let $a, b \in F$. Then by definition $h(a) = a, h(b) = b$ so that $h(a \pm b) = h(a) \pm h(b) = a \pm b, h(ab) = h(a)h(b) = ab$ and $h(a^{-1}) = h(a)^{-1} = a^{-1}$, so that F is closed, hence a subfield of K .

Note that it is not important in this proposition that H actually be a *subgroup* of $\text{Aut}(K)$ — the collection of elements of K fixed by all the elements of a *subset* of $\text{Aut}(K)$ is also a subfield of K .

Definition. If H is a subgroup of the group of automorphisms of K , the subfield of K fixed by all the elements of H is called the *fixed field* of H .

Proposition 4. The association of groups to fields and fields to groups defined above is inclusion reversing, namely

- (1) if $F_1 \subseteq F_2 \subseteq K$ are two subfields of K then $\text{Aut}(K/F_2) \leq \text{Aut}(K/F_1)$, and
- (2) if $H_1 \leq H_2 \leq \text{Aut}(K)$ are two subgroups of automorphisms with associated fixed fields F_1 and F_2 , respectively, then $F_2 \subseteq F_1$.

Proof: Any automorphism of K that fixes F_2 also fixes its subfield F_1 , which gives (1). The second assertion is proved similarly.

Examples

- (1) Suppose $K = \mathbb{Q}(\sqrt{2})$ as in Example 1 above. Then the fixed field of $\text{Aut}(\mathbb{Q}(\sqrt{2})) = \text{Aut}(\mathbb{Q}(\sqrt{2})/\mathbb{Q}) = \{1, \sigma\}$ will be the set of elements of $\mathbb{Q}(\sqrt{2})$ with

$$\sigma(a + b\sqrt{2}) = a + b\sqrt{2}$$

since everything is fixed by the identity automorphism. This is the equation

$$a - b\sqrt{2} = a + b\sqrt{2}.$$

which is equivalent to $b = 0$, so the fixed field of $\text{Aut}(\mathbb{Q}(\sqrt{2})/\mathbb{Q})$ is just \mathbb{Q} .

- (2) Suppose now that $K = \mathbb{Q}(\sqrt[3]{2})$ as in Example 2 above. In this case $\text{Aut}(K) = 1$, so that every element of K is fixed, i.e., the fixed field of $\text{Aut}(\mathbb{Q}(\sqrt[3]{2})/\mathbb{Q})$ is $\mathbb{Q}(\sqrt[3]{2})$.

Given a subfield F of K , the associated group is the collection of automorphisms of K which fix F . Given a group of automorphisms of K , the associated extension is defined by taking F to be the fixed field of the automorphisms. In the first example above, starting with the subfield \mathbb{Q} of $\mathbb{Q}(\sqrt{2})$ one obtains the group $\{1, \sigma\}$ and starting with the group $\{1, \sigma\}$ one obtains the subfield \mathbb{Q} , so there is a “duality” between the two. In the second example, however, starting with the subfield \mathbb{Q} of $\mathbb{Q}(\sqrt[3]{2})$ one obtains only the trivial group and starting with the trivial group one obtains the full field $\mathbb{Q}(\sqrt[3]{2})$.

An examination of the two examples suggests that for the second example there are “not enough” automorphisms to force the fixed field to be \mathbb{Q} rather than the full $\mathbb{Q}(\sqrt[3]{2})$. This in turn seems to be due to the fact that the other roots of $x^3 - 2$, which are the only possible images of $\sqrt[3]{2}$ under an automorphism, are not elements of $\mathbb{Q}(\sqrt[3]{2})$. (Although even if they were we would need to check that the additional maps we could define were *automorphisms*.) We now make precise the notion of fields with “enough” automorphisms (leading to the definition of a *Galois* extension). As one might suspect even from these two examples (and we prove in the next section) these are related to splitting fields.

We first investigate the size of the automorphism group in the case of splitting fields.

Let F be a field and let E be the splitting field over F of $f(x) \in F[x]$. The main tool is Theorem 13.27 on the existence of extensions of isomorphisms, which states that any isomorphism $\varphi : F \xrightarrow{\sim} F'$ of F with F' can be extended to an isomorphism $\sigma : E \xrightarrow{\sim} E'$ between E and the splitting field E' for $f'(x) = \varphi(f(x)) \in F'[x]$.

We now show by induction on $[E : F]$ that the number of such extensions is at most $[E : F]$, with equality if $f(x)$ is separable over F . If $[E : F] = 1$ then $E = F$, $E' = F'$, $\sigma = \varphi$ and the number of extensions is 1. If $[E : F] > 1$ then $f(x)$ has at least one irreducible factor $p(x)$ of degree > 1 with corresponding irreducible factor $p'(x)$ of $f'(x)$. Let α be a fixed root of $p(x)$. If σ is any extension of φ to E , then σ restricted to the subfield $F(\alpha)$ of E is an isomorphism τ of $F(\alpha)$ with some subfield of E' . The isomorphism τ is completely determined by its action on α , i.e., by $\tau\alpha$, since α generates $F(\alpha)$ over F . Just as in Proposition 2, we see that $\tau\alpha$ must be some root β of $p'(\alpha)$. Then we have a diagram

$$\begin{array}{ccccc} \sigma : & E & \xrightarrow{\sim} & E' \\ & | & & | \\ \tau : & F(\alpha) & \xrightarrow{\sim} & F'(\beta) \\ & | & & | \\ \varphi : & F & \xrightarrow{\sim} & F' \end{array}$$

Conversely, for any β a root of $p'(\alpha)$ there are extensions τ and σ giving such a diagram (this is Theorem 13.8 and Theorem 13.27). Hence to count the number of extensions σ we need only count the possible number of these diagrams.

The number of extensions of φ to an isomorphism τ is equal to the number of distinct roots β of $p'(\alpha)$. Since the degree of $p(x)$ and $p'(\alpha)$ are both equal to $[F(\alpha) : F]$, we see that the number of extensions of φ to a τ is at most $[F(\alpha) : F]$, with equality if the roots of $p(x)$ are distinct.

Since E is also the splitting field of $f(x)$ over $F(\alpha)$, E' is the splitting field of $f'(\alpha)$

over $F'(\beta)$, and $[E : F(\alpha)] < [E : F]$, we may apply our induction hypothesis to these field extensions. By induction, the number of extensions of τ to σ is $\leq [E : F(\alpha)]$, with equality if $f(x)$ has distinct roots.

From $[E : F] = [E : F(\alpha)][F(\alpha) : F]$ it follows that the number of extensions of φ to σ is $\leq [E : F]$. We have equality if $p(x)$ and $f(x)$ have distinct roots, which is equivalent to $f(x)$ having distinct roots since $p(x)$ is a factor of $f(x)$, completing the proof by induction.

In the particular case when $F = F'$ and φ is the identity map we have $f(x) = f'(x)$ and $E = E'$ so the isomorphisms of E to E' restricting to φ on F are the automorphisms of E fixing F . We state this as follows:

Proposition 5. Let E be the splitting field over F of the polynomial $f(x) \in F[x]$. Then

$$|\text{Aut}(E/F)| \leq [E : F]$$

with equality if $f(x)$ is separable over F .

Remark: While we were primarily interested in counting the automorphisms of E which fix F (which is the situation of $F = F'$, $\varphi = 1$ above), it would still have been necessary to consider the situation of more general φ (and different fields F') because of the induction step in the proof (which involves the fields $F(\alpha)$ and $F(\beta)$ for two roots of the same polynomial $p(x)$).

One can modify the proof above to show more generally that $|\text{Aut}(K/F)| \leq [K : F]$ for *any* finite extension K/F (we shall prove this in the next section from a slightly different point of view). This gives us a notion of field extensions with “enough” automorphisms.

Definition. Let K/F be a finite extension. Then K is said to be *Galois* over F and K/F is a *Galois* extension if $|\text{Aut}(K/F)| = [K : F]$. If K/F is Galois the group of automorphisms $\text{Aut}(K/F)$ is called the *Galois group* of K/F , denoted $\text{Gal}(K/F)$.

Remark: The Galois group of an extension K/F is sometimes defined to be the group of automorphisms $\text{Aut}(K/F)$ for all K/F . We have chosen the definition above so that the notation $\text{Gal}(K/F)$ will emphasize that the extension K/F has the maximal number of automorphisms.

Corollary 6. If K is the splitting field over F of a separable polynomial $f(x)$ then K/F is Galois.

We shall see in the next section that the converse is also true, which will completely characterize Galois extensions.

Note also that Corollary 6 implies that the splitting field of *any* polynomial over \mathbb{Q} is Galois, since the splitting field of $f(x)$ is clearly the same as the splitting field of the product of the irreducible factors of $f(x)$ (i.e., the polynomial obtained by removing multiple factors), which is separable (Corollary 13.34).

Definition. If $f(x)$ is a separable polynomial over F , then the *Galois group of $f(x)$ over F* is the Galois group of the splitting field of $f(x)$ over F .

Examples

- (1) The extension $\mathbb{Q}(\sqrt{2})/\mathbb{Q}$ is Galois with Galois group $\text{Gal}(\mathbb{Q}(\sqrt{2})/\mathbb{Q}) = \{1, \sigma\} \cong \mathbb{Z}/2\mathbb{Z}$ where σ is the automorphism

$$\begin{aligned}\sigma : \mathbb{Q}(\sqrt{2}) &\xrightarrow{\sim} \mathbb{Q}(\sqrt{2}) \\ a + b\sqrt{2} &\mapsto a - b\sqrt{2}.\end{aligned}$$

- (2) More generally, any quadratic extension K of any field F of characteristic different from 2 is Galois. This follows from the discussion of quadratic extensions following Corollary 13.13, which shows that any extension K of degree 2 of F (where the characteristic of F is not 2) is of the form $F(\sqrt{D})$ for some D hence is the splitting field of $x^2 - D$ (since if $\sqrt{D} \in K$ then also $-\sqrt{D} \in K$).
- (3) The extension $\mathbb{Q}(\sqrt[3]{2})/\mathbb{Q}$ is not Galois since its group of automorphisms is only of order 1.
- (4) The extension $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ is Galois over \mathbb{Q} since it is the splitting field of the polynomial $(x^2 - 2)(x^2 - 3)$. Any automorphism σ is completely determined by its action on the generators $\sqrt{2}$ and $\sqrt{3}$, which must be mapped to $\pm\sqrt{2}$ and $\pm\sqrt{3}$, respectively. Hence the only possibilities for automorphisms are the maps

$$\begin{cases} \sqrt{2} \mapsto \sqrt{2} & \begin{cases} \sqrt{2} \mapsto -\sqrt{2} \\ \sqrt{3} \mapsto \sqrt{3} \end{cases} & \begin{cases} \sqrt{2} \mapsto \sqrt{2} \\ \sqrt{3} \mapsto -\sqrt{3} \end{cases} & \begin{cases} \sqrt{2} \mapsto -\sqrt{2} \\ \sqrt{3} \mapsto -\sqrt{3} \end{cases} \end{cases}.$$

Since the Galois group is of order 4, all these elements are in fact automorphisms of $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ over \mathbb{Q} .

Define the automorphisms σ and τ by

$$\sigma : \begin{cases} \sqrt{2} \mapsto -\sqrt{2} \\ \sqrt{3} \mapsto \sqrt{3} \end{cases} \quad \tau : \begin{cases} \sqrt{2} \mapsto \sqrt{2} \\ \sqrt{3} \mapsto -\sqrt{3} \end{cases}$$

or, more explicitly, by

$$\begin{aligned}\sigma : a + b\sqrt{2} + c\sqrt{3} + d\sqrt{6} &\mapsto a - b\sqrt{2} + c\sqrt{3} - d\sqrt{6} \\ \tau : a + b\sqrt{2} + c\sqrt{3} + d\sqrt{6} &\mapsto a + b\sqrt{2} - c\sqrt{3} - d\sqrt{6}\end{aligned}$$

(since, for example,

$$\sigma(\sqrt{6}) = \sigma(\sqrt{2}\sqrt{3}) = \sigma(\sqrt{2})\sigma(\sqrt{3}) = (-\sqrt{2})(\sqrt{3}) = -\sqrt{6} \quad).$$

Then $\sigma^2(\sqrt{2}) = \sigma(\sigma\sqrt{2}) = \sigma(-\sqrt{2}) = \sqrt{2}$ and clearly $\sigma^2(\sqrt{3}) = \sqrt{3}$. Hence $\sigma^2 = 1$ is the identity automorphism. Similarly, $\tau^2 = 1$. The automorphism $\sigma\tau$ can be easily computed:

$$\sigma\tau(\sqrt{2}) = \sigma(\tau(\sqrt{2})) = \sigma(\sqrt{2}) = -\sqrt{2}$$

and

$$\sigma\tau(\sqrt{3}) = \sigma(\tau(\sqrt{3})) = \sigma(-\sqrt{3}) = -\sqrt{3}$$

so that $\sigma\tau$ is the remaining nontrivial automorphism in the Galois group. Since this automorphism also evidently has order 2 in the Galois group, we have

$$\text{Gal}(\mathbb{Q}(\sqrt{2}, \sqrt{3})/\mathbb{Q}) = \{1, \sigma, \tau, \sigma\tau\}$$

i.e., the Galois group is isomorphic to the Klein 4-group.

Associated to each subgroup of $\text{Gal}(\mathbb{Q}(\sqrt{2}, \sqrt{3})/\mathbb{Q})$ is the corresponding fixed subfield of $\mathbb{Q}(\sqrt{2}, \sqrt{3})$. For example, the subfield corresponding to $\{1, \sigma\tau\}$ is the set of elements fixed by the map

$$\sigma\tau : a + b\sqrt{2} + c\sqrt{3} + d\sqrt{6} \mapsto a - b\sqrt{2} - c\sqrt{3} + d\sqrt{6}$$

which is the set of elements $a + d\sqrt{6}$, i.e., the field $\mathbb{Q}(\sqrt{6})$. One can similarly determine the fixed fields for the other subgroups of the Galois group:

subgroup	fixed field
$\{1\}$	$\mathbb{Q}(\sqrt{2}, \sqrt{3})$
$\{1, \sigma\}$	$\mathbb{Q}(\sqrt{3})$
$\{1, \sigma\tau\}$	$\mathbb{Q}(\sqrt{6})$
$\{1, \tau\}$	$\mathbb{Q}(\sqrt{2})$
$\{1, \sigma, \tau, \sigma\tau\}$	\mathbb{Q}

- (5) The splitting field of $x^3 - 2$ over \mathbb{Q} is Galois of degree 6. The roots of this equation are $\sqrt[3]{2}, \rho \sqrt[3]{2}, \rho^2 \sqrt[3]{2}$ where $\rho = \zeta_3 = \frac{-1 + \sqrt{-3}}{2}$ is a primitive cube root of unity. Hence the splitting field can be written $\mathbb{Q}(\sqrt[3]{2}, \rho \sqrt[3]{2})$. Any automorphism maps each of these two elements to one of the roots of $x^3 - 2$, giving 9 possibilities, but since the Galois group has order 6 not every such map is an automorphism of the field.

To determine the Galois group we use a more convenient set of generators, namely $\sqrt[3]{2}$ and ρ . Then any automorphism σ maps $\sqrt[3]{2}$ to one of $\sqrt[3]{2}, \rho \sqrt[3]{2}, \rho^2 \sqrt[3]{2}$ and maps ρ to ρ or $\rho^2 = \frac{-1 - \sqrt{-3}}{2}$ since these are the roots of the cyclotomic polynomial $\Phi_3(x) = x^2 + x + 1$. Since σ is completely determined by its action on these two elements this gives only 6 possibilities and so each of these possibilities is actually an automorphism. To give these automorphisms explicitly, let σ and τ be the automorphisms defined by

$$\sigma : \begin{cases} \sqrt[3]{2} \mapsto \rho \sqrt[3]{2} \\ \rho \mapsto \rho \end{cases} \quad \tau : \begin{cases} \sqrt[3]{2} \mapsto \sqrt[3]{2} \\ \rho \mapsto \rho^2 = -1 - \rho. \end{cases}$$

As before, these can be given explicitly on the elements of $\mathbb{Q}(\sqrt[3]{2}, \rho)$, which are linear combinations of the basis $\{1, \sqrt[3]{2}, (\sqrt[3]{2})^2, \rho, \rho \sqrt[3]{2}, \rho(\sqrt[3]{2})^2\}$. For example

$$\begin{aligned} \sigma(\rho \sqrt[3]{2}) &= (\rho)(\rho \sqrt[3]{2}) = \rho^2 \sqrt[3]{2} = (-1 - \rho) \sqrt[3]{2} \\ &= -\sqrt[3]{2} - \rho \sqrt[3]{2} \end{aligned}$$

and we may similarly determine the action of σ on the other basis elements. This gives

$$\begin{aligned} \sigma : \quad a + b\sqrt[3]{2} + c\sqrt[3]{4} + d\rho + e\rho \sqrt[3]{2} + f\rho \sqrt[3]{4} &\mapsto \\ a - e\sqrt[3]{2} + (f - c)\sqrt[3]{4} + d\rho + (b - e)\rho \sqrt[3]{2} - c\rho \sqrt[3]{4}. & \end{aligned} \tag{14.1}$$

The other elements of the Galois group are

$$1 : \begin{cases} \sqrt[3]{2} \mapsto \sqrt[3]{2} \\ \rho \mapsto \rho \end{cases} \quad \sigma^2 : \begin{cases} \sqrt[3]{2} \mapsto \rho^2 \sqrt[3]{2} \\ \rho \mapsto \rho \end{cases}$$

Computing $\sigma\tau$ we have

$$\sigma\tau : \begin{cases} \sqrt[3]{2} \xrightarrow{\tau} \sqrt[3]{2} \xrightarrow{\sigma} \rho \sqrt[3]{2} \\ \rho \xrightarrow{\tau} \rho^2 \xrightarrow{\sigma} \rho^2 \end{cases}$$

i.e.,

$$\sigma\tau : \begin{cases} \sqrt[3]{2} \mapsto \rho \sqrt[3]{2} \\ \rho \mapsto \rho^2 \end{cases}$$

so that $\sigma\tau = \tau\sigma^2$. Similarly one computes that $\sigma^3 = \tau^2 = 1$. Hence

$$\text{Gal}(\mathbb{Q}(\sqrt[3]{2}, \zeta_3)/\mathbb{Q}) = \langle \sigma, \tau \rangle \cong S_3$$

is the symmetric group on 3 letters. Alternatively (and less computationally), since $G = \text{Gal}(\mathbb{Q}(\sqrt[3]{2}, \zeta_3)/\mathbb{Q})$ acts as permutations of the 3 roots of $x^3 - 2$, G is a subgroup of S_3 , hence must be S_3 since it is of order 6. The computations above explicitly identify the automorphisms in G and give an explicit isomorphism of G with S_3 .

As in the previous example we can determine the fixed fields for any of the subgroups of the Galois group. For example, consider the fixed field of the subgroup $\{1, \sigma, \sigma^2\}$ generated by σ . These are just the elements fixed by σ (given explicitly in equation (1)) since if an element is fixed by σ then it is also fixed by σ^2 . (In general, the fixed field of some subgroup is the field fixed by a set of generators for the subgroup.) The elements fixed by σ are those with

$$a = a \quad b = -e \quad c = f - c \quad d = d \quad e = b - e \quad f = -c$$

which is equivalent to $b = c = f = e = 0$. Hence the fixed field of $\{1, \sigma, \sigma^2\}$ is the field $\mathbb{Q}(\rho)$.

Remark: This example shows that some care must be exercised in determining Galois groups from the actions on generators. As mentioned, not every map taking $\sqrt[3]{2}$ and $\rho \sqrt[3]{2}$ to roots of $x^3 - 2$ gives rise to an automorphism of the field (for example, the map

$$\begin{aligned} \sqrt[3]{2} &\mapsto \rho \sqrt[3]{2} \\ \rho \sqrt[3]{2} &\mapsto \rho^2 \sqrt[3]{2} \end{aligned}$$

clearly cannot be an automorphism since it is evidently not an injection). The point is that there may be (sometimes very subtle) algebraic relations among the generators and these relations must be respected by an automorphism. For example, the quotient of the generators here is ρ , which is mapped to 1 and not to a root of the minimal polynomial for ρ . Put another way, the quotient of these generators satisfies a quadratic equation and this map does not respect that property.

For another (less trivial) example, compare with the discussion of the splitting field of $x^8 - 2$ in Section 2.

- (6) As in Example 3, the field $\mathbb{Q}(\sqrt[4]{2})$ is not Galois over \mathbb{Q} since any automorphism is determined by where it sends $\sqrt[4]{2}$ and of the four possibilities $\{\pm\sqrt[4]{2}, \pm i\sqrt[4]{2}\}$, only two are elements of the field (the two real roots).

Note that we have

$$\begin{array}{ccccccc} & & & 4 & & & \\ & \overbrace{\quad}^4 & & & & & \\ \mathbb{Q} & \subset & \mathbb{Q}(\sqrt{2}) & \subset & \mathbb{Q}(\sqrt[4]{2}) & \subset & \mathbb{Q}(\sqrt[4]{2}) \\ & \underbrace{\quad}_2 & & & \underbrace{\quad}_2 & & \end{array}$$

where $\mathbb{Q}(\sqrt{2})/\mathbb{Q}$ and $\mathbb{Q}(\sqrt[4]{2})/\mathbb{Q}(\sqrt{2})$ are both Galois extensions by Example 2 since both are quadratic extensions. This shows that a Galois extension of a Galois extension is not necessarily Galois.

- (7) The extension of finite fields $\mathbb{F}_{p^n}/\mathbb{F}_p$ constructed after Proposition 13.37 is Galois by Corollary 6 since \mathbb{F}_{p^n} is the splitting field over \mathbb{F}_p of the separable polynomial $x^{p^n} - x$. It follows that the group of automorphisms for this extension is of order n . The injective homomorphism

$$\begin{aligned} \sigma : \mathbb{F}_{p^n} &\rightarrow \mathbb{F}_{p^n} \\ \alpha &\mapsto \alpha^p \end{aligned}$$

of Proposition 13.35 is surjective in this case since \mathbb{F}_{p^n} is finite, hence is an isomorphism. This gives an automorphism of \mathbb{F}_{p^n} , called the *Frobenius* automorphism, which we shall denote by σ_p . Iterating σ_p we have $\sigma_p^2(\alpha) = \sigma_p(\sigma_p(\alpha)) = (\alpha^p)^p = \alpha^{p^2}$. Similarly we have

$$\sigma_p^i(\alpha) = \alpha^{p^i} \quad i = 0, 1, 2, \dots$$

Since $\alpha^{p^n} = \alpha$, we see that $\sigma_p^{p^n} = 1$ is the identity automorphism. No lower power of σ_p can be the identity, since this would imply $\alpha^{p^i} = \alpha$ for all $\alpha \in \mathbb{F}_{p^n}$ for some $i < n$, which is impossible since there are only p^i roots of this equation. It follows that σ_p is of order n in the Galois group, which means that $\text{Gal}(\mathbb{F}_{p^n}/\mathbb{F}_p)$ is *cyclic* of order n , with the Frobenius automorphism σ_p as generator.

- (8) The inseparable extension $\mathbb{F}_2(x)$ over $\mathbb{F}_2(t)$ where $x^2 - t = 0$ considered in Section 13.5 is not Galois. Any automorphism of this degree 2 extension is determined by its action on x , which must be sent to a root of the equation $x^2 - t$. We have already seen that there is only one root of this equation (with multiplicity 2) since we are in a field of characteristic 2. Hence the extension has only the trivial automorphism. Note that $\mathbb{F}_2(x)$ is the splitting field for $x^2 - t$ over $\mathbb{F}_2(t)$, so this example shows the separability condition in Corollary 6 is necessary.

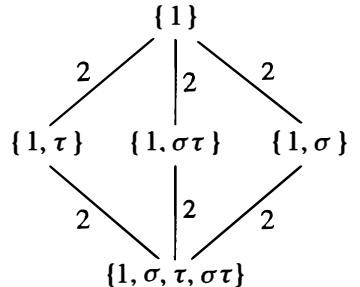
EXERCISES

- (a) Show that if the field K is generated over F by the elements $\alpha_1, \dots, \alpha_n$ then an automorphism σ of K fixing F is uniquely determined by $\sigma(\alpha_1), \dots, \sigma(\alpha_n)$. In particular show that an automorphism fixes K if and only if it fixes a set of generators for K .
- (b) Let $G \leq \text{Gal}(K/F)$ be a subgroup of the Galois group of the extension K/F and suppose $\sigma_1, \dots, \sigma_k$ are generators for G . Show that the subfield E/F is fixed by G if and only if it is fixed by the generators $\sigma_1, \dots, \sigma_k$.

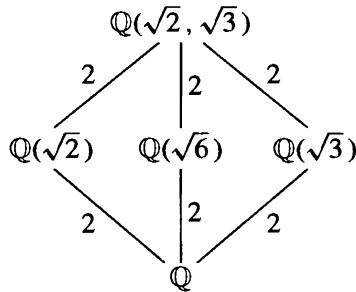
2. Let τ be the map $\tau : \mathbb{C} \rightarrow \mathbb{C}$ defined by $\tau(a + bi) = a - bi$ (*complex conjugation*). Prove that τ is an automorphism of \mathbb{C} .
3. Determine the fixed field of complex conjugation on \mathbb{C} .
4. Prove that $\mathbb{Q}(\sqrt{2})$ and $\mathbb{Q}(\sqrt{3})$ are not isomorphic.
5. Determine the automorphisms of the extension $\mathbb{Q}(\sqrt[4]{2})/\mathbb{Q}(\sqrt{2})$ explicitly.
6. Let k be a field.
 - (a) Show that the mapping $\varphi : k[t] \rightarrow k[t]$ defined by $\varphi(f(t)) = f(at + b)$ for fixed $a, b \in k$, $a \neq 0$ is an automorphism of $k[t]$ which is the identity on k .
 - (b) Conversely, let φ be an automorphism of $k[t]$ which is the identity on k . Prove that there exist $a, b \in k$ with $a \neq 0$ such that $\varphi(f(t)) = f(at + b)$ as in (a).
7. This exercise determines $\text{Aut}(\mathbb{R}/\mathbb{Q})$.
 - (a) Prove that any $\sigma \in \text{Aut}(\mathbb{R}/\mathbb{Q})$ takes squares to squares and takes positive reals to positive reals. Conclude that $a < b$ implies $\sigma a < \sigma b$ for every $a, b \in \mathbb{R}$.
 - (b) Prove that $-\frac{1}{m} < a - b < \frac{1}{m}$ implies $-\frac{1}{m} < \sigma a - \sigma b < \frac{1}{m}$ for every positive integer m . Conclude that σ is a continuous map on \mathbb{R} .
 - (c) Prove that any continuous map on \mathbb{R} which is the identity on \mathbb{Q} is the identity map, hence $\text{Aut}(\mathbb{R}/\mathbb{Q}) = 1$.
8. Prove that the automorphisms of the rational function field $k(t)$ which fix k are precisely the *fractional linear transformations* determined by $t \mapsto \frac{at + b}{ct + d}$ for $a, b, c, d \in k$, $ad - bc \neq 0$ (so $f(t) \in k(t)$ maps to $f(\frac{at + b}{ct + d})$) (cf. Exercise 18 of Section 13.2).
9. Determine the fixed field of the automorphism $t \mapsto t + 1$ of $k(t)$.
10. Let K be an extension of the field F . Let $\varphi : K \rightarrow K'$ be an isomorphism of K with a field K' which maps F to the subfield F' of K' . Prove that the map $\sigma \mapsto \varphi\sigma\varphi^{-1}$ defines a group isomorphism $\text{Aut}(K/F) \xrightarrow{\sim} \text{Aut}(K'/F')$.

14.2 THE FUNDAMENTAL THEOREM OF GALOIS THEORY

In the Galois extension $\text{Gal}(\mathbb{Q}(\sqrt{2}, \sqrt{3})/\mathbb{Q})$ considered in the previous section, there was a strong similarity between the diagram of subgroups of the Galois group:



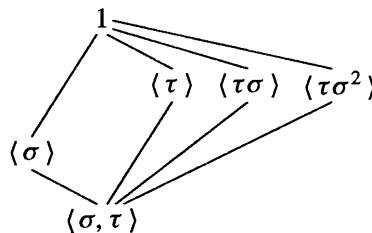
and the diagram of corresponding fixed fields



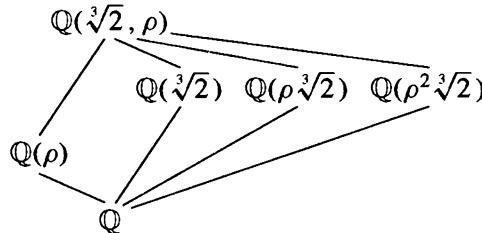
(we have inverted the lattice of subgroups because of the inclusion-reversing nature of the correspondence).

Note that this is also the diagram of *all* known subfields of the extension and that in this case each of the subfields is also a Galois extension of \mathbb{Q} .

In a similar way there is a strong similarity between the diagram



of subgroups of the Galois group and the diagram of known subfields for the splitting field of $x^3 - 2$:



where the subfields in the second diagram are precisely the fixed fields of the subgroups in the first diagram.

Note in this pair of diagrams only the subgroup $\langle \sigma \rangle$ generated by σ is normal in S_3 and that the subfield $\mathbb{Q}(\rho)$ is the only subfield Galois over \mathbb{Q} .

The Fundamental Theorem of Galois Theory states that the relations observed in the two examples above are not coincidental and hold for any Galois extension. Before proving this we first develop some preliminary results on *group characters*, of which field automorphisms give particular examples.

Definition. A *character*¹ χ of a group G with values in a field L is a homomorphism from G to the multiplicative group of L :

$$\chi : G \rightarrow L^\times$$

i.e., $\chi(g_1g_2) = \chi(g_1)\chi(g_2)$ for all $g_1, g_2 \in G$ and $\chi(g)$ is a nonzero element of L for all $g \in G$.

Definition. The characters $\chi_1, \chi_2, \dots, \chi_n$ of G are said to be *linearly independent* over L if they are linearly independent as functions on G , i.e., if there is no nontrivial relation

$$a_1\chi_1 + a_2\chi_2 + \cdots + a_n\chi_n = 0 \quad (a_1, \dots, a_n \in L \text{ not all } 0) \quad (14.2)$$

as a function on G (that is, $a_1\chi_1(g) + a_2\chi_2(g) + \cdots + a_n\chi_n(g) = 0$ for all $g \in G$).

Theorem 7. (Linear Independence of Characters) If $\chi_1, \chi_2, \dots, \chi_n$ are distinct characters of G with values in L then they are linearly independent over L .

Proof: Suppose the characters were linearly dependent. Among all the linear dependence relations (2) above, choose one with the minimal number m of nonzero coefficients a_i . We may suppose (by renumbering, if necessary) that the m nonzero coefficients are a_1, a_2, \dots, a_m :

$$a_1\chi_1 + a_2\chi_2 + \cdots + a_m\chi_m = 0.$$

Then for any $g \in G$ we have

$$a_1\chi_1(g) + a_2\chi_2(g) + \cdots + a_m\chi_m(g) = 0. \quad (14.3)$$

Let g_0 be an element with $\chi_1(g_0) \neq \chi_m(g_0)$ (which exists, since $\chi_1 \neq \chi_m$). Since (3) holds for every element of G , in particular we have

$$a_1\chi_1(g_0g) + a_2\chi_2(g_0g) + \cdots + a_m\chi_m(g_0g) = 0$$

i.e.,

$$a_1\chi_1(g_0)\chi_1(g) + a_2\chi_2(g_0)\chi_2(g) + \cdots + a_m\chi_m(g_0)\chi_m(g) = 0. \quad (14.4)$$

Multiplying equation (3) by $\chi_m(g_0)$ and subtracting from equation (4) we obtain

$$\begin{aligned} [\chi_m(g_0) - \chi_1(g_0)]a_1\chi_1(g) + [\chi_m(g_0) - \chi_2(g_0)]a_2\chi_2(g) + \cdots \\ + [\chi_m(g_0) - \chi_{m-1}(g_0)]a_{m-1}\chi_{m-1}(g) = 0, \end{aligned}$$

which holds for all $g \in G$. But the first coefficient is nonzero and this is a relation with fewer nonzero coefficients, a contradiction.

Consider now an injective homomorphism σ of a field K into a field L , called an *embedding* of K into L . Then in particular σ is a homomorphism of the multiplicative group $G = K^\times$ into the multiplicative group L^\times , so σ may be viewed as a character of K^\times with values in L . Note also that this character contains all of the useful information about the values of σ viewed simply as a *function* on K , since the only point of K not considered in K^\times is 0, and we know σ maps 0 to 0.

¹This is the definition of a *linear* character. More general characters will be studied in Chapter 18.

Corollary 8. If $\sigma_1, \sigma_2, \dots, \sigma_n$ are distinct embeddings of a field K into a field L , then they are linearly independent as functions on K . In particular distinct automorphisms of a field K are linearly independent as functions on K .

We now use Corollary 8 to prove the fundamental relation between the orders of subgroups of the automorphism group of a field K and the degrees of the extensions over their fixed fields.

Theorem 9. Let $G = \{\sigma_1 = 1, \sigma_2, \dots, \sigma_n\}$ be a subgroup of automorphisms of a field K and let F be the fixed field. Then

$$[K : F] = n = |G|.$$

Proof: Suppose first that $n > [K : F]$ and let $\omega_1, \omega_2, \dots, \omega_m$ be a basis for K over F ($m = [K : F]$). Then the system

$$\sigma_1(\omega_1)x_1 + \sigma_2(\omega_1)x_2 + \cdots + \sigma_n(\omega_1)x_n = 0$$

⋮

$$\sigma_1(\omega_m)x_1 + \sigma_2(\omega_m)x_2 + \cdots + \sigma_n(\omega_m)x_n = 0$$

of m equations in n unknowns x_1, x_2, \dots, x_n has a nontrivial solution $\beta_1, \beta_2, \dots, \beta_n$ in K since by assumption there are more unknowns than equations.

Let a_1, a_2, \dots, a_m be m arbitrary elements of F . The field F is by definition fixed by $\sigma_1, \dots, \sigma_n$ so each of these elements is fixed by every σ_i , i.e., $\sigma_i(a_j) = a_j$, $i = 1, 2, \dots, n$, $j = 1, 2, \dots, m$. Multiplying the first equation above by a_1 , the second by a_2, \dots , the last by a_m then gives the system of equations

$$\sigma_1(a_1\omega_1)\beta_1 + \sigma_2(a_1\omega_1)\beta_2 + \cdots + \sigma_n(a_1\omega_1)\beta_n = 0$$

⋮

$$\sigma_1(a_m\omega_m)\beta_1 + \sigma_2(a_m\omega_m)\beta_2 + \cdots + \sigma_n(a_m\omega_m)\beta_n = 0.$$

Adding these equations we see that there are elements β_1, \dots, β_n in K , not all 0, satisfying

$$\sigma_1(a_1\omega_1 + a_2\omega_2 + \cdots + a_m\omega_m)\beta_1 + \cdots + \sigma_n(a_1\omega_1 + a_2\omega_2 + \cdots + a_m\omega_m)\beta_n = 0$$

for all choices of a_1, \dots, a_m in F . Since $\omega_1, \dots, \omega_m$ is an F -basis for K , every $\alpha \in K$ is of the form $a_1\omega_1 + a_2\omega_2 + \cdots + a_m\omega_m$, so the previous equation means

$$\sigma_1(\alpha)\beta_1 + \cdots + \sigma_n(\alpha)\beta_n = 0$$

for all $\alpha \in K$. But this means the distinct automorphisms $\sigma_1, \dots, \sigma_n$ are linearly dependent over K , contradicting Corollary 8.

We have proved $n \leq [K : F]$. Note that we have so far not used the fact that $\sigma_1, \sigma_2, \dots, \sigma_n$ are the elements of a group.

Suppose now that $n < [K : F]$. Then there are more than n F -linearly independent elements of K , say $\alpha_1, \dots, \alpha_{n+1}$. The system

$$\sigma_1(\alpha_1)x_1 + \sigma_1(\alpha_2)x_2 + \cdots + \sigma_1(\alpha_{n+1})x_{n+1} = 0$$

⋮

$$\sigma_n(\alpha_1)x_1 + \sigma_n(\alpha_2)x_2 + \cdots + \sigma_n(\alpha_{n+1})x_{n+1} = 0$$

(14.5)

of n equations in $n + 1$ unknowns x_1, \dots, x_{n+1} has a solution $\beta_1, \dots, \beta_{n+1}$ in K where not all the β_i , $i = 1, 2, \dots, n + 1$ are 0. If all the elements of the solution $\beta_1, \dots, \beta_{n+1}$ were elements of F then the first equation (recall $\sigma_1 = 1$ is the identity automorphism) would contradict the linear independence over F of $\alpha_1, \alpha_2, \dots, \alpha_{n+1}$. Hence at least one β_i , $i = 1, 2, \dots, n + 1$, is not an element of F .

Among all the nontrivial solutions $(\beta_1, \dots, \beta_{n+1})$ of the system (5) choose one with the minimal number r of nonzero β_i . By renumbering if necessary we may assume β_1, \dots, β_r are nonzero. Dividing the equations by β_r we may also assume $\beta_r = 1$. We have already seen that at least one of $\beta_1, \dots, \beta_{r-1}, 1$ is not an element of F (which shows in particular that $r > 1$), say $\beta_1 \notin F$. Then our system of equations reads

$$\begin{aligned} \sigma_1(\alpha_1)\beta_1 + \cdots + \sigma_1(\alpha_{r-1})\beta_{r-1} + \sigma_1(\alpha_r) &= 0 \\ &\vdots \\ \sigma_n(\alpha_1)\beta_1 + \cdots + \sigma_n(\alpha_{r-1})\beta_{r-1} + \sigma_n(\alpha_r) &= 0 \end{aligned} \tag{14.6}$$

or more briefly

$$\sigma_i(\alpha_1)\beta_1 + \cdots + \sigma_i(\alpha_{r-1})\beta_{r-1} + \sigma_i(\alpha_r) = 0 \quad i = 1, 2, \dots, n. \tag{14.7}$$

Since $\beta_1 \notin F$, there is an automorphism σ_{k_0} ($k_0 \in \{1, 2, \dots, n\}$) with $\sigma_{k_0}\beta_1 \neq \beta_1$. If we apply the automorphism σ_{k_0} to the equations in (6), we obtain the system of equations

$$\sigma_{k_0}\sigma_j(\alpha_1)\sigma_{k_0}(\beta_1) + \cdots + \sigma_{k_0}\sigma_j(\alpha_{r-1})\sigma_{k_0}(\beta_{r-1}) + \sigma_{k_0}\sigma_j(\alpha_r) = 0 \tag{14.8}$$

for $j = 1, 2, \dots, n$. But the elements

$$\sigma_{k_0}\sigma_1, \sigma_{k_0}\sigma_2, \dots, \sigma_{k_0}\sigma_n$$

are the same as the elements

$$\sigma_1, \sigma_2, \dots, \sigma_n$$

in some order since these elements form a *group*. In other words, if we define the index i by $\sigma_{k_0}\sigma_j = \sigma_i$ then i and j both run over the set $\{1, 2, \dots, n\}$. Hence the equations in (8) can be written

$$\sigma_i(\alpha_1)\sigma_{k_0}(\beta_1) + \cdots + \sigma_i(\alpha_{r-1})\sigma_{k_0}(\beta_{r-1}) + \sigma_i(\alpha_r) = 0. \tag{14.8'}$$

If we now subtract the equations in (8') from those in (7) we obtain the system

$$\sigma_i(\alpha_1)[\beta_1 - \sigma_{k_0}(\beta_1)] + \cdots + \sigma_i(\alpha_{r-1})[\beta_{r-1} - \sigma_{k_0}(\beta_{r-1})] = 0$$

for $i = 1, 2, \dots, n$. But this is a solution to the system of equations (5) with

$$x_1 = \beta_1 - \sigma_{k_0}(\beta_1) \neq 0$$

(by the choice of k_0), hence is nontrivial and has fewer than r nonzero x_i . This is a contradiction and completes the proof.

Our first use of this result is to prove that the inequality of Proposition 5 holds for any finite extension K/F .

Corollary 10. Let K/F be any finite extension. Then

$$|\text{Aut}(K/F)| \leq [K : F]$$

with equality if and only if F is the fixed field of $\text{Aut}(K/F)$. Put another way, K/F is Galois if and only if F is the fixed field of $\text{Aut}(K/F)$.

Proof: Let F_1 be the fixed field of $\text{Aut}(K/F)$, so that

$$F \subseteq F_1 \subseteq K.$$

By Theorem 9, $[K : F_1] = |\text{Aut}(K/F)|$. Hence $[K : F] = |\text{Aut}(K/F)|[F_1 : F]$, which proves the corollary.

Corollary 11. Let G be a finite subgroup of automorphisms of a field K and let F be the fixed field. Then every automorphism of K fixing F is contained in G , i.e., $\text{Aut}(K/F) = G$, so that K/F is Galois, with Galois group G .

Proof: By definition F is fixed by all the elements of G so we have $G \leq \text{Aut}(K/F)$ (and the question is whether there are any automorphisms of K fixing F not in G i.e., whether this containment is proper). Hence $|G| \leq |\text{Aut}(K/F)|$. By the theorem we have $|G| = [K : F]$ and by the previous corollary $|\text{Aut}(K/F)| \leq [K : F]$. This gives

$$[K : F] = |G| \leq |\text{Aut}(K/F)| \leq [K : F]$$

and it follows that we must have equalities throughout, proving the corollary.

Corollary 12. If $G_1 \neq G_2$ are distinct finite subgroups of automorphisms of a field K then their fixed fields are also distinct.

Proof: Suppose F_1 is the fixed field of G_1 and F_2 is the fixed field of G_2 . If $F_1 = F_2$ then by definition F_1 is fixed by G_2 . By the previous corollary any automorphism fixing F_1 is contained in G_1 , hence $G_2 \leq G_1$. Similarly $G_1 \leq G_2$ and so $G_1 = G_2$.

By the corollaries above we see that taking the fixed fields for distinct finite subgroups of $\text{Aut}(K)$ gives distinct subfields of K over which K is Galois. Further, the degrees of the extensions are given by the orders of the subgroups. We saw this explicitly for the fields $K = \mathbb{Q}(\sqrt{2}, \sqrt{3})$ and $K = \mathbb{Q}(\sqrt[3]{2}, \rho)$ above. A portion of the Fundamental Theorem states that these are *all* the subfields of K .

The next result provides the converse of Proposition 5 and characterizes Galois extensions.

Theorem 13. The extension K/F is Galois if and only if K is the splitting field of some separable polynomial over F . Furthermore, if this is the case then every irreducible polynomial with coefficients in F which has a root in K is separable and has all its roots in K (so in particular K/F is a separable extension).

Proof: Proposition 5 proves that the splitting field of a separable polynomial is Galois.

We now show that if K/F is Galois then every irreducible polynomial $p(x)$ in $F[x]$ having a root in K splits completely in K . Set $G = \text{Gal}(K/F)$. Let $\alpha \in K$ be a root of $p(x)$ and consider the elements

$$\alpha, \sigma_2(\alpha), \dots, \sigma_n(\alpha) \in K \quad (14.9)$$

where $\{1, \sigma_2, \dots, \sigma_n\}$ are the elements of $\text{Gal}(K/F)$. Let

$$\alpha, \alpha_2, \alpha_3, \dots, \alpha_r$$

denote the *distinct* elements in (9). If $\tau \in G$ then since G is a group the elements $\{\tau, \tau\sigma_2, \dots, \tau\sigma_n\}$ are the same as the elements $\{1, \sigma_2, \dots, \sigma_n\}$ in some order. It follows that applying $\tau \in G$ to the elements in (9) simply permutes them, so in particular applying τ to $\alpha, \alpha_2, \alpha_3, \dots, \alpha_r$ also permutes these elements. The polynomial

$$f(x) = (x - \alpha)(x - \alpha_2) \cdots (x - \alpha_r)$$

therefore has coefficients which are fixed by all the elements of G since the elements of G simply permute the factors. Hence the coefficients lie in the fixed field of G , which by Corollary 10 is the field F . Hence $f(x) \in F[x]$.

Since $p(x)$ is irreducible and has α as a root, $p(x)$ is the minimal polynomial for α over F , hence divides any polynomial with coefficients in F having α as a root (this is Proposition 13.9). It follows that $p(x)$ divides $f(x)$ in $F[x]$ and since $f(x)$ obviously divides $p(x)$ in $K[x]$ by Proposition 2, we have

$$p(x) = f(x).$$

In particular, this shows that $p(x)$ is separable and that all its roots lie in K (in fact they are among the elements $\alpha, \sigma_2\alpha, \dots, \sigma_n\alpha$), proving the last statement of the theorem.

To complete the proof, suppose K/F is Galois and let $\omega_1, \omega_2, \dots, \omega_n$ be a basis for K/F . Let $p_i(x)$ be the minimal polynomial for ω_i over F , $i = 1, 2, \dots, n$. Then by what we have just proved, $p_i(x)$ is separable and has all its roots in K . Let $g(x)$ be the polynomial obtained by removing any multiple factors in the product $p_1(x) \cdots p_n(x)$ (the “squarefree part”). Then the splitting field of the two polynomials is the same and this field is K (all the roots lie in K , so K contains the splitting field, but $\omega_1, \omega_2, \dots, \omega_n$ are among the roots, so the splitting field contains K). Hence K is the splitting field of the separable polynomial $g(x)$.

Definition. Let K/F be a Galois extension. If $\alpha \in K$ the elements $\sigma\alpha$ for σ in $\text{Gal}(K/F)$ are called the *conjugates* (or *Galois conjugates*) of α over F . If E is a subfield of K containing F , the field $\sigma(E)$ is called the *conjugate field* of E over F .

The proof of the theorem shows that in a Galois extension K/F the other roots of the minimal polynomial over F of any element $\alpha \in K$ are precisely the distinct conjugates of α under the Galois group of K/F .

The second statement in this theorem also shows that K is not Galois over F if we can find even one irreducible polynomial over F having a root in K but not having *all* its roots in K . This justifies in a very strong sense the intuition from earlier examples that Galois extensions are extensions with “enough” distinct roots of irreducible polynomials (namely, if it contains one root then it contains all the roots).

- Finally, notice that we now have 4 characterizations of Galois extensions K/F :
- (1) splitting fields of separable polynomials over F
 - (2) fields where F is precisely the set of elements fixed by $\text{Aut}(K/F)$ (in general, the fixed field may be larger than F)
 - (3) fields with $[K : F] = |\text{Aut}(K/F)|$ (the original definition)
 - (4) finite, normal and separable extensions.

Theorem 14. (Fundamental Theorem of Galois Theory) Let K/F be a Galois extension and set $G = \text{Gal}(K/F)$. Then there is a bijection

$$\left\{ \begin{array}{c} \text{subfields } E \\ \text{of } K \\ \text{containing } F \end{array} \middle| \begin{array}{c} K \\ E \\ F \end{array} \right\} \leftrightarrow \left\{ \begin{array}{c} \text{subgroups } H \\ \text{of } G \\ \mid \\ G \end{array} \middle| \begin{array}{c} 1 \\ H \\ \mid \end{array} \right\}$$

given by the correspondences

$$\begin{array}{ccc} E & \rightarrow & \left\{ \begin{array}{c} \text{the elements of } G \\ \text{fixing } E \end{array} \right\} \\ \left\{ \begin{array}{c} \text{the fixed field} \\ \text{of } H \end{array} \right\} & \leftarrow & H \end{array}$$

which are inverse to each other. Under this correspondence,

- (1) (inclusion reversing) If E_1, E_2 correspond to H_1, H_2 , respectively, then $E_1 \subseteq E_2$ if and only if $H_2 \leq H_1$
- (2) $[K : E] = |H|$ and $[E : F] = |G : H|$, the index of H in G :

$$\begin{array}{ccc} K & & |H| \\ | & \} & \\ E & & \\ | & \} & |G : H| \\ F & & \end{array}$$

- (3) K/E is always Galois, with Galois group $\text{Gal}(K/E) = H$:

$$\begin{array}{ccc} K & & \\ | & H & \\ E & & \end{array}$$

- (4) E is Galois over F if and only if H is a normal subgroup in G . If this is the case, then the Galois group is isomorphic to the quotient group

$$\text{Gal}(E/F) \cong G/H.$$

More generally, even if H is not necessarily normal in G , the isomorphisms of E (into a fixed algebraic closure of F containing K) which fix F are in one to one correspondence with the cosets $\{\sigma H\}$ of H in G .

- (5) If E_1, E_2 correspond to H_1, H_2 , respectively, then the intersection $E_1 \cap E_2$ corresponds to the group $\langle H_1, H_2 \rangle$ generated by H_1 and H_2 and the composite field $E_1 E_2$ corresponds to the intersection $H_1 \cap H_2$. Hence the lattice of subfields

of K containing F and the lattice of subgroups of G are “dual” (the lattice diagram for one is the lattice diagram for the other turned upside down).

Proof: Given any subgroup H of G we obtain a unique fixed field $E = K_H$ by Corollary 12. This shows that the correspondence above is injective from right to left.

If K is the splitting field of the separable polynomial $f(x) \in F[x]$ then we may also view $f(x)$ as an element of $E[x]$ for any subfield E of K containing F . Then K is also the splitting field of $f(x)$ over E , so the extension K/E is Galois. By Corollary 10, E is the fixed field of $\text{Aut}(K/E) \leq G$, showing that *every* subfield of K containing F arises as the fixed field for some subgroup of G . Hence the correspondence above is surjective from right to left, hence a bijection. The correspondences are inverse to each other since the automorphisms fixing E are precisely $\text{Aut}(K/E)$ by Corollary 10.

We have already seen that the Galois correspondence is inclusion reversing in Proposition 4, which gives (1).

If $E = K_H$ is the fixed field of H , then Theorem 9 gives $[K : E] = |H|$ and $[K : F] = |G|$. Taking the quotient gives $[E : F] = |G : H|$, which proves (2).

Corollary 11 gives (3) immediately.

Suppose $E = K_H$ is the fixed field of the subgroup H . Every $\sigma \in G = \text{Gal}(K/F)$ when restricted to E is an embedding $\sigma|_E$ of E with the subfield $\sigma(E)$ of K . Conversely, let $\tau : E \xrightarrow{\sim} \overline{F}$ be any embedding of E (into a fixed algebraic closure \overline{F} of F containing K) which fixes F . Then $\tau(E)$ is in fact contained in K : if $\alpha \in E$ has minimal polynomial $m_\alpha(x)$ over F then $\tau(\alpha)$ is another root of $m_\alpha(x)$ and K contains all these roots by Theorem 13. As above K is the splitting field of $f(x)$ over E and so also the splitting field of $\tau f(x)$ (which is the same as $f(x)$ since $f(x)$ has coefficients in F) over $\tau(E)$. Theorem 13.27 on extending isomorphisms then shows that we can extend τ to an isomorphism σ :

$$\begin{array}{ccc} \sigma : & K & \xrightarrow{\sim} & K \\ & | & & | \\ \tau : & E & \xrightarrow{\sim} & \tau(E). \end{array}$$

Since σ fixes F (because τ does), it follows that *every* embedding τ of E fixing F is the restriction to E of some automorphism σ of K fixing F , in other words, every embedding of E is of the form $\sigma|_E$ for some $\sigma \in G$.

Two automorphisms $\sigma, \sigma' \in G$ restrict to the *same* embedding of E if and only if $\sigma^{-1}\sigma'$ is the identity map on E . But then $\sigma^{-1}\sigma' \in H$ (i.e., $\sigma' \in \sigma H$) since by (3) the automorphisms of K which fix E are precisely the elements in H . Hence the distinct embeddings of E are in bijection with the cosets σH of H in G . In particular this gives

$$|\text{Emb}(E/F)| = [G : H] = [E : F]$$

where $\text{Emb}(E/F)$ denotes the set of embeddings of E (into a fixed algebraic closure of F) which fix F . Note that $\text{Emb}(E/F)$ contains the automorphisms $\text{Aut}(E/F)$.

The extension E/F will be Galois if and only if $|\text{Aut}(E/F)| = [E : F]$. By the equality above, this will be the case if and only if each of the *embeddings* of E is actually an *automorphism* of E , i.e., if and only if $\sigma(E) = E$ for every $\sigma \in G$.

If $\sigma \in G$, then the subgroup of G fixing the field $\sigma(E)$ is the group $\sigma H \sigma^{-1}$, i.e.,

$$\sigma(E) = K_{\sigma H \sigma^{-1}}.$$

To see this observe that if $\sigma\alpha \in \sigma(E)$ then

$$(\sigma h\sigma^{-1})(\sigma\alpha) = \sigma(h\alpha) = \sigma\alpha \quad \text{for all } h \in H,$$

since h fixes $\alpha \in E$, which shows that $\sigma H\sigma^{-1}$ fixes $\sigma(E)$. The group fixing $\sigma(E)$ has order equal to the degree of K over $\sigma(E)$. But this is the same as the degree of K over E since the fields are isomorphic, hence the same as the order of H . Hence $\sigma H\sigma^{-1}$ is precisely the group fixing $\sigma(E)$ since we have shown containment and their orders are the same.

Because of the bijective nature of the Galois correspondence already proved we know that two subfields of K containing F are equal if and only if their fixing subgroups are equal in G . Hence $\sigma(E) = E$ for all $\sigma \in G$ if and only if $\sigma H\sigma^{-1} = H$ for all $\sigma \in G$, in other words E is Galois over F if and only if H is a normal subgroup of G .

We have already identified the embeddings of E over F as the set of cosets of H in G and when H is normal in G seen that the embeddings are automorphisms. It follows that in this case the *group* of cosets G/H is identified with the *group* of automorphisms of the Galois extension E/F by the definition of the group operation (composition of automorphisms). Hence $G/H \cong \text{Gal}(E/F)$ when H is normal in G , which completes the proof of (4).

Suppose H_1 is the subgroup of elements of G fixing the subfield E_1 and H_2 is the subgroup of elements of G fixing the subfield E_2 . Any element in $H_1 \cap H_2$ fixes both E_1 and E_2 , hence fixes every element in the composite E_1E_2 , since the elements in this field are algebraic combinations of the elements of E_1 and E_2 . Conversely, if an automorphism σ fixes the composite E_1E_2 then in particular σ fixes E_1 , i.e., $\sigma \in H_1$, and σ fixes E_2 , i.e., $\sigma \in H_2$, hence $\sigma \in H_1 \cap H_2$. This proves that the composite E_1E_2 corresponds to the intersection $H_1 \cap H_2$. Similarly, the intersection $E_1 \cap E_2$ corresponds to the group (H_1, H_2) generated by H_1 and H_2 , completing the proof of the theorem.

Example: $(\mathbb{Q}(\sqrt{2}, \sqrt{3})$ and $\mathbb{Q}(\sqrt[3]{2}, \rho))$

We have already seen examples of this theorem at the beginning of this section. We now see that the diagrams of subfields for the two fields $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ and $\mathbb{Q}(\sqrt[3]{2}, \rho)$ given before indicate *all* the subfields for these two fields.

Since every subgroup of the Klein 4-group is normal, all the subfields of $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ are Galois extensions of \mathbb{Q} .

Similarly, since the only nontrivial normal subgroup of S_3 is the subgroup of order 3, we see that only the subfield $\mathbb{Q}(\rho)$ of $K = \mathbb{Q}(\sqrt[3]{2}, \rho)$ is Galois over \mathbb{Q} , with Galois group isomorphic to $S_3/(\sigma)$, i.e., the cyclic group of order 2. For example, the nontrivial automorphism of $\mathbb{Q}(\rho)$ is induced by restricting any element (τ , for instance) in the nontrivial coset of (σ) to $\mathbb{Q}(\rho)$. This is clear from the explicit descriptions of these automorphisms given before — each of the elements τ , $\tau\sigma$, $\tau\sigma^2$ in this coset map ρ to ρ^2 . The restrictions of the elements of $\text{Gal}(K/\mathbb{Q})$ to the (non-Galois) cubic subfields do not give automorphisms of these fields in general, rather giving isomorphisms of these fields with each other, in accordance with (4) of the theorem.

Example: $(\mathbb{Q}(\sqrt{2} + \sqrt{3}))$

Consider the field $\mathbb{Q}(\sqrt{2} + \sqrt{3})$. This is clearly a subfield of the Galois extension $\mathbb{Q}(\sqrt{2}, \sqrt{3})$. The other roots of the minimal polynomial for $\sqrt{2} + \sqrt{3}$ over \mathbb{Q} are therefore

the distinct conjugates of $\sqrt{2} + \sqrt{3}$ under the Galois group. The conjugates are

$$\pm\sqrt{2} \pm \sqrt{3}$$

which are easily seen to be distinct. The minimal polynomial is therefore

$$[x - (\sqrt{2} + \sqrt{3})][x - (\sqrt{2} - \sqrt{3})][x - (-\sqrt{2} + \sqrt{3})][x - (-\sqrt{2} - \sqrt{3})]$$

which is quickly computed to be the polynomial $x^4 - 10x^2 + 1$. It follows that this polynomial is irreducible and that

$$\mathbb{Q}(\sqrt{2}, \sqrt{3}) = \mathbb{Q}(\sqrt{2} + \sqrt{3}),$$

either by degree considerations or by noting that only the automorphism 1 of $\{1, \sigma, \tau, \sigma\tau\}$ fixes $\sqrt{2} + \sqrt{3}$ so the fixing group for this field is the same as for $\mathbb{Q}(\sqrt{2}, \sqrt{3})$.

Example: (Splitting Field of $x^8 - 2$)

The splitting field of $x^8 - 2$ over \mathbb{Q} is generated by $\theta = \sqrt[8]{2}$ (any fixed 8th root of 2, say the real one) and a primitive 8th root of unity $\zeta = \zeta_8$. Recall from Section 13.6 that

$$\mathbb{Q}(\zeta_8) = \mathbb{Q}(i, \sqrt{2}).$$

Since $\theta^4 = \sqrt{2}$ we see that the splitting field is generated by θ and i . The subfield $\mathbb{Q}(\theta)$ is of degree 8 over \mathbb{Q} (since $x^8 - 2$ is irreducible, being Eisenstein), and all the elements of this field are real. Hence $i \notin \mathbb{Q}(\theta)$ and since i generates at most a quadratic extension of this field, the splitting field

$$\mathbb{Q}(\sqrt[8]{2}, \zeta_8) = \mathbb{Q}(\sqrt[8]{2}, i)$$

is of degree 16 over \mathbb{Q} .

The Galois group is determined by the action on the generators θ and i which gives the possibilities

$$\begin{cases} \theta \mapsto \zeta^a \theta & a = 0, 1, 2, \dots, 7 \\ i \mapsto \pm i \end{cases}$$

Since we have already seen that the degree of the extension is 16 and there are only 16 possible such maps, it follows that in fact each of the maps above is an automorphism of $\mathbb{Q}(\sqrt[8]{2}, i)$ over \mathbb{Q} .

Define the two automorphisms

$$\sigma : \begin{cases} \theta \mapsto \zeta \theta \\ i \mapsto i \end{cases} \quad \tau : \begin{cases} \theta \mapsto \theta \\ i \mapsto -i \end{cases}$$

(τ is the map induced by complex conjugation). Since

$$\begin{aligned} \zeta = \zeta_8 &= \frac{\sqrt{2}}{2} + i \frac{\sqrt{2}}{2} = \frac{1}{2}(1+i)\sqrt{2} \\ &= \frac{1}{2}(1+i)\theta^4 \end{aligned}$$

we can easily compute what happens to ζ from the explicit expressions for the powers of ζ in the following Figure 1.

Using these explicit values we find

$$\sigma : \begin{cases} \theta \mapsto \zeta \theta \\ i \mapsto i \\ \zeta \mapsto -\zeta = \zeta^5 \end{cases} \quad \tau : \begin{cases} \theta \mapsto \theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^7 \end{cases}$$

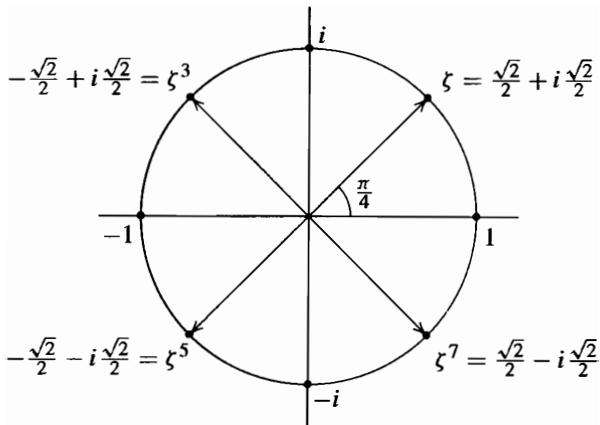


Fig. 1

Note that the reason we are interested in also keeping track of the action on the element ζ is that it will be needed in computing the composites of automorphisms, for example in computing

$$\begin{aligned}\sigma^2(\theta) &= \sigma(\zeta\theta) = \sigma(\zeta)\sigma(\theta) = (-\zeta)(\zeta\theta) = -\zeta^2\theta \\ &= -i\theta.\end{aligned}$$

We can similarly compute the following automorphisms:

$$\begin{array}{lll} \sigma : \begin{cases} \theta \mapsto \zeta\theta \\ i \mapsto i \\ \zeta \mapsto \zeta^5 \end{cases} & \tau\sigma : \begin{cases} \theta \mapsto \zeta^7\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^3 \end{cases} \\ \sigma^2 : \begin{cases} \theta \mapsto \zeta^6\theta \\ i \mapsto i \\ \zeta \mapsto \zeta \end{cases} & \tau\sigma^2 : \begin{cases} \theta \mapsto \zeta^2\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^7 \end{cases} \\ \sigma^3 : \begin{cases} \theta \mapsto \zeta^7\theta \\ i \mapsto i \\ \zeta \mapsto -\zeta \end{cases} & \tau\sigma^3 : \begin{cases} \theta \mapsto \zeta\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^3 \end{cases} \\ \sigma^4 : \begin{cases} \theta \mapsto -\theta \\ i \mapsto i \\ \zeta \mapsto \zeta \end{cases} & \tau\sigma^4 : \begin{cases} \theta \mapsto -\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^7 \end{cases} \\ \sigma^5 : \begin{cases} \theta \mapsto \zeta^5\theta \\ i \mapsto i \\ \zeta \mapsto -\zeta \end{cases} & \tau\sigma^5 : \begin{cases} \theta \mapsto \zeta^3\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^3 \end{cases} \\ \sigma^6 : \begin{cases} \theta \mapsto \zeta^2\theta \\ i \mapsto i \\ \zeta \mapsto \zeta \end{cases} & \tau\sigma^6 : \begin{cases} \theta \mapsto \zeta^6\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^7 \end{cases} \end{array}$$

$$\sigma^7 : \begin{cases} \theta \mapsto \zeta^3\theta \\ i \mapsto i \\ \zeta \mapsto -\zeta \end{cases} \quad \tau\sigma^7 : \begin{cases} \theta \mapsto \zeta^5\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^3. \end{cases}$$

Since this exhausts the possibilities, these elements (together with 1 and τ) are the Galois group. We see in particular that σ and τ generate the Galois group. To determine the relations satisfied by these elements, we observe first that clearly $\tau^2 = 1$ and $(\sigma^4)^2 = 1$, so that

$$\sigma^8 = \tau^2 = 1.$$

Also, we compute

$$\sigma\tau : \begin{cases} \theta \mapsto \zeta\theta \\ i \mapsto -i \\ \zeta \mapsto \zeta^3 \end{cases}$$

so that

$$\sigma\tau = \tau\sigma^3.$$

It is not too difficult to show that these relations define the group completely, i.e.,

$$\text{Gal}(\mathbb{Q}(\sqrt[8]{2}, i)/\mathbb{Q}) = \langle \sigma, \tau \mid \sigma^8 = \tau^2 = 1, \sigma\tau = \tau\sigma^3 \rangle.$$

Such a group is called a *quasidihedral group* (recall that the dihedral group of order 16 would have the relation $\sigma\tau = \tau\sigma^7$ instead of $\sigma\tau = \tau\sigma^3$) and is a subgroup of S_8 since the Galois group is a subgroup of the permutations of the 8 roots of $x^8 - 2$.

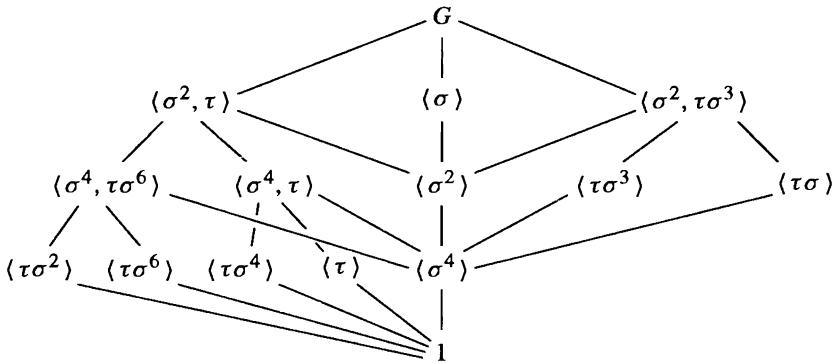
This example again illustrates that one must take care in determining Galois groups from the actions on generators. We first computed the degree of the Galois extension above to determine the number of elements in the Galois group. Had we proceeded directly from the original generators $\theta = \sqrt[8]{2}$ and $\zeta = \zeta_8$ we might have (incorrectly) concluded that there were a total of 32 elements in the Galois group, since the first generator is mapped to any of 8 possible roots of $x^8 - 2$ and the second generator is mapped to any of 4 possible roots of its minimal polynomial $\Phi_4(x) = x^4 + 1$. The problem, as previously indicated, is that these choices are not independent. Here the reason is provided by the algebraic relation

$$\theta^4 = \sqrt{2} = \zeta + \zeta^7$$

which shows that one cannot specify the images of θ and ζ independently — their images must again satisfy this algebraic relation. This relation is perhaps sufficiently subtle to serve as a caution against rashly concluding maps are automorphisms. We note that in general it is necessary to provide justification that maps are automorphisms. This can be accomplished for example by using the extension theorems or by using degree considerations as we did here.

Determining the lattice of subgroups of this group G is a straightforward problem.

The lattice is the following:



Determining the subfields corresponding to these subgroups (which by the Fundamental Theorem gives *all* the subfields of $\mathbb{Q}(\sqrt[8]{2}, i)$) is quite simple for a number of the subgroups above using (2) of the Fundamental Theorem, which states that the degree of the extension over \mathbb{Q} is equal to the *index* of the fixing subgroup. It then suffices to find a subfield of the right degree which is fixed by the subgroup in question. Remember also that if a subfield is fixed by the *generators* of a subgroup, then it is fixed by the subgroup. For example, from the explicit description for the automorphism σ we see that $\mathbb{Q}(i)$ is fixed by the group generated by σ . Since this is a subgroup of index 2 and $\mathbb{Q}(i)$ is of degree 2 over \mathbb{Q} , it must be the full fixed field. Most of the fixed fields for the subgroups above can be determined in as simple a manner.

For the subgroups of order 4 on the right (namely, generated by $\tau\sigma^3$ and by $\tau\sigma$), it is perhaps not so easy to see how to determine the corresponding fixed field. For the subgroup H generated by $\tau\sigma^3$ we may proceed as follows: the element $\theta^2 = \sqrt[4]{2}$ is clearly fixed by σ^4 . By the diagram above, σ^4 is a normal subgroup of H of index 2, with representatives 1, $\tau\sigma^3$ for the cosets. Consider the element

$$\alpha = (1 + \tau\sigma^3)\theta^2 = \theta^2 + \tau\sigma^3\theta^2.$$

Then α is fixed by σ^4 (we are in a commutative group H of order 4, so σ^4 commutes with 1 and $\tau\sigma^3$ and we already know θ^2 is fixed by σ^4). But (and this is the point), α is also fixed by $\tau\sigma^3$:

$$\begin{aligned}\tau\sigma^3\alpha &= \tau\sigma^3(1 + \tau\sigma^3)\theta^2 = [\tau\sigma^3 + (\tau\sigma^3)^2]\theta^2 \\ &= (\tau\sigma^3 + \sigma^4)\theta^2\end{aligned}$$

and the last expression is just α since $\sigma^4\theta^2 = \theta^2$. Hence α is an element of the fixed field for H . Explicitly

$$\alpha = \sqrt[4]{2} + i\sqrt[4]{2} = (1 + i)\sqrt[4]{2}.$$

A quick check shows that α is not fixed by the automorphism σ^2 , so by the diagram of subgroups above, it follows that the fixing subgroup for the field $\mathbb{Q}(\alpha)$ is no larger than H , hence is precisely H , which gives us our fixed field. This also gives the fixed field for $\langle \tau\sigma \rangle$ by recalling that in general if E is the fixed field of H then the fixed field of $\tau H \tau^{-1}$ is the field $\tau(E)$. For $H = \langle \tau\sigma^3 \rangle$, $\tau H \tau^{-1} = \langle \tau\sigma \rangle$, with fixed field given by $\tau(\alpha) = (1 - i)\sqrt[4]{2}$.

In general one tries to determine elements which are fixed by a given subgroup H of the Galois group (cf. the exercises, which indicate where the element above arose) and

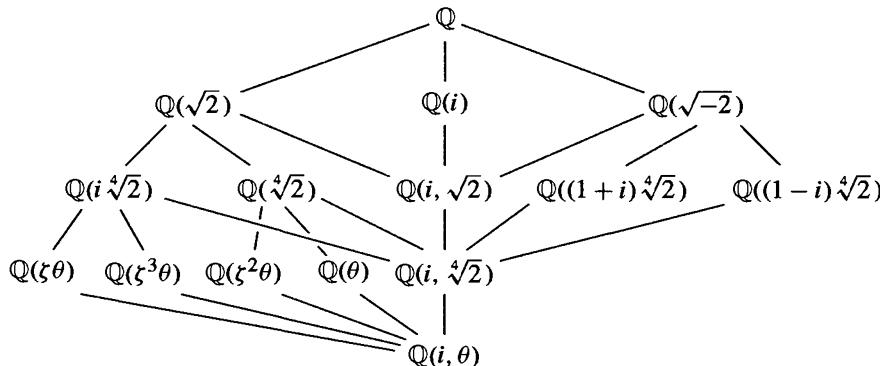
attempts to generate a sufficiently large field to give the full fixed field. In our case we were able to accomplish this with a single generator. We shall see later that every finite extension of \mathbb{Q} is a simple extension, so there will be a single generator of this type, but in general it may be difficult to produce it directly.

The element α is a root of the polynomial

$$x^4 + 8$$

which must therefore be irreducible since we have already determined that a root of this polynomial generates an extension of degree 4 over \mathbb{Q} .

In a similar way it is possible to complete the diagram of subfields of $\mathbb{Q}(\sqrt[8]{2}, i)$, which we have inverted to emphasize its relation with the subgroup diagram above ($\theta = \sqrt[8]{2}$):



Note that the group $\langle \sigma^4 \rangle$ is normal in G (in fact it is the center of G) with quotient $G/\langle \sigma^4 \rangle \cong D_8$, so the corresponding fixed field $\mathbb{Q}(i, \sqrt[4]{2})$ is Galois over \mathbb{Q} with D_8 as Galois group. Being Galois it is a splitting field, evidently the splitting field for $x^4 - 2$. The lattice of subfields for this field is then immediate from the lattice above.

We end this example with the following amusing aspect of this Galois extension. It is an easy exercise to verify that

$$\langle \sigma^2, \tau \rangle \cong D_8 \quad \langle \sigma \rangle \cong \mathbb{Z}/8\mathbb{Z} \quad \langle \sigma^2, \tau\sigma^3 \rangle \cong Q_8$$

where D_8 is the dihedral group of order 8 and Q_8 is the quaternion group of order 8. It follows that the field $\mathbb{Q}(\sqrt[8]{2}, i)$ is Galois of degree 8 over its three quadratic subfields

$$\mathbb{Q}(\sqrt[8]{2}) \quad \mathbb{Q}(i) \quad \mathbb{Q}(\sqrt{-2})$$

with dihedral, cyclic and quaternion Galois groups, respectively, so that three of the 5 possible groups of order 8 (and both non-abelian ones) appear as Galois groups in this extension.

We shall consider additional examples and applications in the following sections.

EXERCISES

1. Determine the minimal polynomial over \mathbb{Q} for the element $\sqrt{2} + \sqrt{5}$.
2. Determine the minimal polynomial over \mathbb{Q} for the element $1 + \sqrt[3]{2} + \sqrt[3]{4}$.
3. Determine the Galois group of $(x^2 - 2)(x^2 - 3)(x^2 - 5)$. Determine *all* the subfields of the splitting field of this polynomial.

4. Let p be a prime. Determine the elements of the Galois group of $x^p - 2$.
5. Prove that the Galois group of $x^p - 2$ for p a prime is isomorphic to the group of matrices $\begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}$ where $a, b \in \mathbb{F}_p$, $a \neq 0$.
6. Let $K = \mathbb{Q}(\sqrt[8]{2}, i)$ and let $F_1 = \mathbb{Q}(i)$, $F_2 = \mathbb{Q}(\sqrt{2})$, $F_3 = \mathbb{Q}(\sqrt{-2})$. Prove that $\text{Gal}(K/F_1) \cong \mathbb{Z}_8$, $\text{Gal}(K/F_2) \cong D_8$, $\text{Gal}(K/F_3) \cong Q_8$.
7. Determine all the subfields of the splitting field of $x^8 - 2$ which are Galois over \mathbb{Q} .
8. Suppose K is a Galois extension of F of degree p^n for some prime p and some $n \geq 1$. Show there are Galois extensions of F contained in K of degrees p and p^{n-1} .
9. Give an example of fields F_1, F_2, F_3 with $\mathbb{Q} \subset F_1 \subset F_2 \subset F_3$, $[F_3 : \mathbb{Q}] = 8$ and each field is Galois over all its subfields with the exception that F_2 is not Galois over \mathbb{Q} .
10. Determine the Galois group of the splitting field over \mathbb{Q} of $x^8 - 3$.
11. Suppose $f(x) \in \mathbb{Z}[x]$ is an irreducible quartic whose splitting field has Galois group S_4 over \mathbb{Q} (there are many such quartics, cf. Section 6). Let θ be a root of $f(x)$ and set $K = \mathbb{Q}(\theta)$. Prove that K is an extension of \mathbb{Q} of degree 4 which has no proper subfields. Are there any Galois extensions of \mathbb{Q} of degree 4 with no proper subfields?
12. Determine the Galois group of the splitting field over \mathbb{Q} of $x^4 - 14x^2 + 9$.
13. Prove that if the Galois group of the splitting field of a cubic over \mathbb{Q} is the cyclic group of order 3 then all the roots of the cubic are real.
14. Show that $\mathbb{Q}(\sqrt{2 + \sqrt{2}})$ is a cyclic quartic field, i.e., is a Galois extension of degree 4 with cyclic Galois group.
15. (*Biquadratic Extensions*) Let F be a field of characteristic $\neq 2$.
 - (a) If $K = F(\sqrt{D_1}, \sqrt{D_2})$ where $D_1, D_2 \in F$ have the property that none of D_1 , D_2 or $D_1 D_2$ is a square in F , prove that K/F is a Galois extension with $\text{Gal}(K/F)$ isomorphic to the Klein 4-group.
 - (b) Conversely, suppose K/F is a Galois extension with $\text{Gal}(K/F)$ isomorphic to the Klein 4-group. Prove that $K = F(\sqrt{D_1}, \sqrt{D_2})$ where $D_1, D_2 \in F$ have the property that none of D_1 , D_2 or $D_1 D_2$ is a square in F .
16. (a) Prove that $x^4 - 2x^2 - 2$ is irreducible over \mathbb{Q} .
 (b) Show the roots of this quartic are

$$\alpha_1 = \sqrt{1 + \sqrt{3}} \quad \alpha_3 = -\sqrt{1 + \sqrt{3}}$$

$$\alpha_2 = \sqrt{1 - \sqrt{3}} \quad \alpha_4 = -\sqrt{1 - \sqrt{3}}.$$
 - (c) Let $K_1 = \mathbb{Q}(\alpha_1)$ and $K_2 = \mathbb{Q}(\alpha_2)$. Show that $K_1 \neq K_2$, and $K_1 \cap K_2 = \mathbb{Q}(\sqrt{3}) = F$.
 - (d) Prove that K_1, K_2 and $K_1 K_2$ are Galois over F with $\text{Gal}(K_1 K_2/F)$ the Klein 4-group. Write out the elements of $\text{Gal}(K_1 K_2/F)$ explicitly. Determine all the subgroups of the Galois group and give their corresponding fixed subfields of $K_1 K_2$ containing F .
 - (e) Prove that the splitting field of $x^4 - 2x^2 - 2$ over \mathbb{Q} is of degree 8 with dihedral Galois group.

The following two exercises indicate one method for constructing elements in subfields of a given field and are quite useful in many computations.

17. Let K/F be any finite extension and let $\alpha \in K$. Let L be a Galois extension of F containing K and let $H \leq \text{Gal}(L/F)$ be the subgroup corresponding to K . Define the *norm* of α from

K to F to be

$$N_{K/F}(\alpha) = \prod_{\sigma} \sigma(\alpha),$$

where the product is taken over all the embeddings of K into an algebraic closure of F (so over a set of coset representatives for H in $\text{Gal}(L/F)$ by the Fundamental Theorem of Galois Theory). This is a product of Galois conjugates of α . In particular, if K/F is Galois this is $\prod_{\sigma \in \text{Gal}(K/F)} \sigma(\alpha)$.

- (a) Prove that $N_{K/F}(\alpha) \in F$.
- (b) Prove that $N_{K/F}(\alpha\beta) = N_{K/F}(\alpha)N_{K/F}(\beta)$, so that the norm is a multiplicative map from K to F .
- (c) Let $K = F(\sqrt{D})$ be a quadratic extension of F . Show that $N_{K/F}(a + b\sqrt{D}) = a^2 - Db^2$.
- (d) Let $m_{\alpha}(x) = x^d + a_{d-1}x^{d-1} + \cdots + a_1x + a_0 \in F[x]$ be the minimal polynomial for $\alpha \in K$ over F . Let $n = [K : F]$. Prove that d divides n , that there are d distinct Galois conjugates of α which are all repeated n/d times in the product above and conclude that $N_{K/F}(\alpha) = (-1)^n a_0^{n/d}$.

18. With notation as in the previous problem, define the *trace* of α from K to F to be

$$\text{Tr}_{K/F}(\alpha) = \sum_{\sigma} \sigma(\alpha),$$

a sum of Galois conjugates of α .

- (a) Prove that $\text{Tr}_{K/F}(\alpha) \in F$.
- (b) Prove that $\text{Tr}_{K/F}(\alpha + \beta) = \text{Tr}_{K/F}(\alpha) + \text{Tr}_{K/F}(\beta)$, so that the trace is an additive map from K to F .
- (c) Let $K = F(\sqrt{D})$ be a quadratic extension of F . Show that $\text{Tr}_{K/F}(a + b\sqrt{D}) = 2a$.
- (d) Let $m_{\alpha}(x)$ be as in the previous problem. Prove that $\text{Tr}_{K/F}(\alpha) = -\frac{n}{d}a_{d-1}$.

19. With notation as in the previous problems show that $N_{K/F}(a\alpha) = a^n N_{K/F}(\alpha)$ and $\text{Tr}_{K/F}(a\alpha) = a\text{Tr}_{K/F}(\alpha)$ for all a in the base field F . In particular show that $N_{K/F}(a) = a^n$ and $\text{Tr}_{K/F}(a) = na$ for all $a \in F$.

20. With notation as in the previous problems show more generally that $\prod_{\sigma} (x - \sigma(\alpha)) = (m_{\alpha}(x))^{n/d}$.

21. Use the linear independence of characters to show that for any Galois extension K of F there is an element $\alpha \in K$ with $\text{Tr}_{K/F}(\alpha) \neq 0$.

22. Suppose K/F is a Galois extension and let σ be an element of the Galois group.

- (a) Suppose $\alpha \in K$ is of the form $\alpha = \frac{\beta}{\sigma\beta}$ for some nonzero $\beta \in K$. Prove that $N_{K/F}(\alpha) = 1$.
- (b) Suppose $\alpha \in K$ is of the form $\alpha = \beta - \sigma\beta$ for some $\beta \in K$. Prove that $\text{Tr}_{K/F}(\alpha) = 0$.

The next exercise and Exercise 26 following establish the multiplicative and additive forms of Hilbert's Theorem 90. These are instances of the vanishing of a first cohomology group, as will be discussed in Section 17.3.

23. (*Hilbert's Theorem 90*) Let K be a Galois extension of F with cyclic Galois group of order n generated by σ . Suppose $\alpha \in K$ has $N_{K/F}(\alpha) = 1$. Prove that α is of the form $\alpha = \frac{\beta}{\sigma\beta}$ for some nonzero $\beta \in K$. [By the linear independence of characters show there exists some $\theta \in K$ such that

$$\beta = \theta + \alpha\sigma(\theta) + (\alpha\sigma\alpha)\sigma^2(\theta) + \cdots + (\alpha\sigma\alpha\cdots\sigma^{n-2}\alpha)\sigma^{n-1}(\theta)$$

is nonzero. Compute $\frac{\beta}{\sigma\beta}$ using the fact that α has norm 1 to F .]

24. Prove that the rational solutions $a, b \in \mathbb{Q}$ of Pythagoras' equation $a^2 + b^2 = 1$ are of the form $a = \frac{s^2 - t^2}{s^2 + t^2}$ and $b = \frac{2st}{s^2 + t^2}$ for some $s, t \in \mathbb{Q}$ and hence show that any right triangle with integer sides has sides of lengths $(m^2 - n^2, 2mn, m^2 + n^2)$ for some integers m, n . [Note that $a^2 + b^2 = 1$ is equivalent to $N_{\mathbb{Q}(i)/\mathbb{Q}}(a + ib) = 1$, then use Hilbert's Theorem 90 above with $\beta = s + it$.]
25. Generalize the previous problem to determine all the rational solutions of the equation $a^2 + Db^2 = 1$ for $D \in \mathbb{Z}$, $D > 0$, D not a perfect square in \mathbb{Z} .
26. (*Additive Hilbert's Theorem 90*) Let K be a Galois extension of F with cyclic Galois group of order n generated by σ . Suppose $\alpha \in K$ has $\text{Tr}_{K/F}(\alpha) = 0$. Prove that α is of the form $\alpha = \beta - \sigma\beta$ for some $\beta \in K$. [Let $\theta \in K$ be an element with $\text{Tr}_{K/F}(\theta) \neq 0$ by a previous exercise, let

$$\beta = \frac{1}{\text{Tr}_{K/F}(\theta)} [\alpha\sigma(\theta) + (\alpha + \sigma\alpha)\sigma^2(\theta) + \cdots + (\alpha + \sigma\alpha + \cdots + \sigma^{n-2}\alpha)\sigma^{n-1}(\theta)]$$

and compute $\beta - \sigma\beta$.]

27. Let $\alpha = \sqrt{(2 + \sqrt{2})(3 + \sqrt{3})}$ (positive real square roots for concreteness) and consider the extension $E = \mathbb{Q}(\alpha)$.

- (a) Show that $a = (2 + \sqrt{2})(3 + \sqrt{3})$ is not a square in $F = \mathbb{Q}(\sqrt{2}, \sqrt{3})$. [If $a = c^2$, $c \in F$, then $a\varphi(a) = (2 + \sqrt{2})^2(6) = (c\varphi c)^2$ for the automorphism $\varphi \in \text{Gal}(F/\mathbb{Q})$ fixing $\mathbb{Q}(\sqrt{2})$. Since $c\varphi c = N_{F/\mathbb{Q}(\sqrt{2})}(c) \in \mathbb{Q}(\sqrt{2})$ conclude that this implies $\sqrt{6} \in \mathbb{Q}(\sqrt{2})$, a contradiction.]
- (b) Conclude from (a) that $[E : \mathbb{Q}] = 8$. Prove that the roots of the minimal polynomial over \mathbb{Q} for α are the 8 elements $\pm\sqrt{(2 \pm \sqrt{2})(3 \pm \sqrt{3})}$.
- (c) Let $\beta = \sqrt{(2 - \sqrt{2})(3 + \sqrt{3})}$. Show that $\alpha\beta = \sqrt{2}(3 + \sqrt{3}) \in F$ so that $\beta \in E$. Show similarly that the other roots are also elements of E so that E is a Galois extension of \mathbb{Q} . Show that the elements of the Galois group are precisely the maps determined by mapping α to one of the eight elements in (b).
- (d) Let $\sigma \in \text{Gal}(E/\mathbb{Q})$ be the automorphism which maps α to β . Show that since $\sigma(\alpha^2) = \beta^2$ that $\sigma(\sqrt{2}) = -\sqrt{2}$ and $\sigma(\sqrt{3}) = \sqrt{3}$. From $\alpha\beta = \sqrt{2}(3 + \sqrt{3})$ conclude that $\sigma(\alpha\beta) = -\alpha\beta$ and hence $\sigma(\beta) = -\alpha$. Show that σ is an element of order 4 in $\text{Gal}(E/\mathbb{Q})$.
- (e) Show similarly that the map τ defined by $\tau(\alpha) = \sqrt{(2 + \sqrt{2})(3 - \sqrt{3})}$ is an element of order 4 in $\text{Gal}(E/\mathbb{Q})$. Prove that σ and τ generate the Galois group, $\sigma^4 = \tau^4 = 1$, $\sigma^2 = \tau^2$ and that $\sigma\tau = \tau\sigma^3$.
- (f) Conclude that $\text{Gal}(E/\mathbb{Q}) \cong Q_8$, the quaternion group of order 8.
28. Let $f(x) \in F[x]$ be an irreducible polynomial of degree n over the field F , let L be the splitting field of $f(x)$ over F and let α be a root of $f(x)$ in L . If K is any Galois extension of F contained in L , show that the polynomial $f(x)$ splits into a product of m irreducible polynomials each of degree d over K , where $m = [F(\alpha) \cap K : F]$ and $d = [K(\alpha) : K]$ (cf. also the generalization in Exercise 4 of Section 4). [If H is the subgroup of the Galois group of L over F corresponding to K then the factors of $f(x)$ over K correspond to the orbits of H on the roots of $f(x)$. Then use Exercise 9 of Section 4.1.]

29. Let k be a field and let $k(t)$ be the field of rational functions in the variable t . Define the maps σ and τ of $k(t)$ to itself by $\sigma f(t) = f(\frac{1}{1-t})$ and $\tau f(t) = f(\frac{1}{t})$ for $f(t) \in k(t)$.

(a) Prove that σ and τ are automorphisms of $k(t)$ (cf. Exercise 8 of Section 1) and that the group $G = \langle \sigma, \tau \rangle$ they generate is isomorphic to S_3 .

(b) Prove that the element $t = \frac{(t^2 - t + 1)^3}{t^2(t-1)^2}$ is fixed by all the elements of G .

(c) Prove that $k(t)$ is precisely the fixed field of G in $k(t)$ [compute the degree of the extension].

30. Prove that the fixed field of the subgroup of automorphisms generated by τ in the previous problem is $k(t + \frac{1}{t})$. Prove that the fixed field of the subgroup generated by the automorphism $\tau\sigma^2$ (which maps t to $1-t$) is $k(t(1-t))$. Determine the fixed field of the subgroup generated by $\tau\sigma$ and the fixed field of the subgroup generated by σ .

31. Let K be a finite extension of F of degree n . Let α be an element of K .

(a) Prove that α acting by left multiplication on K is an F -linear transformation T_α of K .

(b) Prove that the minimal polynomial for α over F is the same as the minimal polynomial for the linear transformation T_α .

(c) Prove that the trace $\text{Tr}_{K/F}(\alpha)$ is the trace of the $n \times n$ matrix defined by T_α (which justifies these two uses of the same word “trace”). Prove that the norm $N_{K/F}(\alpha)$ is the determinant of T_α .

14.3 FINITE FIELDS

A finite field \mathbb{F} has characteristic p for some prime p so is a finite dimensional vector space over \mathbb{F}_p . If the dimension is n , i.e., $[\mathbb{F} : \mathbb{F}_p] = n$, then \mathbb{F} has precisely p^n elements. We have already seen (following Proposition 13.37) that \mathbb{F} is then isomorphic to the splitting field of the polynomial $x^{p^n} - x$, hence is unique up to isomorphism. We denote the finite field of order p^n by \mathbb{F}_{p^n} .

The field \mathbb{F}_{p^n} is Galois over \mathbb{F}_p , with cyclic Galois group of order n generated by the Frobenius automorphism

$$\text{Gal}(\mathbb{F}_{p^n}/\mathbb{F}_p) = \langle \sigma_p \rangle \cong \mathbb{Z}/n\mathbb{Z}$$

where

$$\begin{aligned} \sigma_p : \mathbb{F}_{p^n} &\rightarrow \mathbb{F}_{p^n} \\ \alpha &\mapsto \alpha^p \end{aligned}$$

(Example 7 following Corollary 6). By the Fundamental Theorem, every subfield of \mathbb{F}_{p^n} corresponds to a subgroup of $\mathbb{Z}/n\mathbb{Z}$. Hence for every divisor d of n there is precisely one subfield of \mathbb{F}_{p^n} of degree d over \mathbb{F}_p , namely the fixed field of the subgroup generated by σ_p^d of order n/d , and there are no other subfields. This field is isomorphic to \mathbb{F}_{p^d} , the unique finite field of order p^d .

Since the Galois group is abelian, every subgroup is normal, so each of the subfields \mathbb{F}_{p^d} (d a divisor of n) is Galois over \mathbb{F}_p (which is also clear from the fact that these are themselves splitting fields). Further, the Galois group $\text{Gal}(\mathbb{F}_{p^d}/\mathbb{F}_p)$ is generated by the image of σ_p in the quotient group $\text{Gal}(\mathbb{F}_{p^n}/\mathbb{F}_p)/\langle \sigma_p^d \rangle$. If we denote this element

again by σ_p , we recover the Frobenius automorphism for the extension $\mathbb{F}_{p^d}/\mathbb{F}_p$. (Note, however, that σ_p has order n in $\text{Gal}(\mathbb{F}_{p^n}/\mathbb{F}_p)$ and order d in $\text{Gal}(\mathbb{F}_{p^d}/\mathbb{F}_p)$.)

We summarize this in the following proposition.

Proposition 15. Any finite field is isomorphic to \mathbb{F}_{p^n} for some prime p and some integer $n \geq 1$. The field \mathbb{F}_{p^n} is the splitting field over \mathbb{F}_p of the polynomial $x^{p^n} - x$, with cyclic Galois group of order n generated by the Frobenius automorphism σ_p . The subfields of \mathbb{F}_{p^n} are all Galois over \mathbb{F}_p and are in one to one correspondence with the divisors d of n . They are the fields \mathbb{F}_{p^d} , the fixed fields of σ_p^d .

The corresponding statements for the finite extensions of any finite field are easy consequences of Proposition 15 and are outlined in the exercises.

As an elementary application we have the following result on the polynomial $x^4 + 1$ in $\mathbb{Z}[x]$.

Corollary 16. The irreducible polynomial $x^4 + 1 \in \mathbb{Z}[x]$ is reducible modulo every prime p .

Proof: Consider the polynomial $x^4 + 1$ over $\mathbb{F}_p[x]$ for the prime p . If $p = 2$ we have $x^4 + 1 = (x + 1)^4$ and the polynomial is reducible. Assume now that p is odd. Then $p^2 - 1$ is divisible by 8 since p is congruent mod 8 to 1, 3, 5 or 7 and all of these square to 1 mod 8. Hence $x^{p^2-1} - 1$ is divisible by $x^8 - 1$. Then we have the divisibilities

$$x^4 + 1 \mid x^8 - 1 \mid x^{p^2-1} - 1 \mid x^{p^2} - x$$

which shows that all the roots of $x^4 + 1$ are roots of $x^{p^2} - x$. (Equivalently, these roots are fixed by the square of the Frobenius automorphism σ_p^2 .) Since the roots of $x^{p^2} - x$ are the field \mathbb{F}_{p^2} , it follows that the extension generated by any root of $x^4 + 1$ is at most of degree 2 over \mathbb{F}_p , which means that $x^4 + 1$ cannot be irreducible over \mathbb{F}_p .

The multiplicative group $\mathbb{F}_{p^n}^\times$ is obviously a finite subgroup of the multiplicative group of a field. By Proposition 9.18, this is a *cyclic* group. If θ is any generator, then clearly $\mathbb{F}_{p^n} = \mathbb{F}_p(\theta)$. This proves the following result.

Proposition 17. The finite field \mathbb{F}_{p^n} is simple. In particular, there exists an irreducible polynomial of degree n over \mathbb{F}_p for every $n \geq 1$.

We have described the finite fields \mathbb{F}_{p^n} above as the splitting fields of the polynomials $x^{p^n} - x$. By the previous proposition, this field can also be described as a quotient of $\mathbb{F}_p[x]$, namely by the minimal polynomial for θ . Since θ is necessarily a root of $x^{p^n} - x$, we see that the minimal polynomial for θ is a divisor of $x^{p^n} - x$ of degree n .

Conversely, let $p(x)$ be any irreducible polynomial of degree d , say, dividing $x^{p^n} - x$. If α is a root of $p(x)$, then the extension $\mathbb{F}_p(\alpha)$ is a subfield of \mathbb{F}_{p^n} of degree d . Hence d is a divisor of n and the extension is Galois by Proposition 15 (in fact, the extension \mathbb{F}_{p^d}) so in particular all the roots of $p(x)$ are contained in $\mathbb{F}_p(\alpha)$.

The elements of \mathbb{F}_{p^n} are precisely the roots of $x^{p^n} - x$. If we group together the factors $x - \alpha$ of this polynomial according to the degree d of their minimal polynomials over \mathbb{F}_p , we obtain

Proposition 18. The polynomial $x^{p^n} - x$ is precisely the product of all the distinct irreducible polynomials in $\mathbb{F}_p[x]$ of degree d where d runs through all divisors of n .

This proposition can be used to produce irreducible polynomials over \mathbb{F}_p recursively. For example, the irreducible quadratics over \mathbb{F}_2 are the divisors of

$$\frac{x^4 - x}{x(x-1)}$$

which gives the single polynomial $x^2 + x + 1$. Similarly, the irreducible cubics over this field are the divisors of

$$\frac{x^8 - x}{x(x-1)} = x^6 + x^5 + x^4 + x^3 + x^2 + x + 1$$

which factors into the two cubics $x^3 + x + 1$ and $x^3 + x^2 + 1$. The irreducible quartics are given by dividing $x^{16} - x$ by $x(x-1)$ and the irreducible quadratic $x^2 + x + 1$ above and then factoring into irreducible quartics:

$$\frac{x^{16} - x}{x(x-1)(x^2 + x + 1)} = (x^4 + x^3 + x^2 + x + 1)(x^4 + x^3 + 1)(x^4 + x + 1).$$

This gives a method for determining the product of all the irreducible polynomials over \mathbb{F}_p of a given degree. There exist efficient algorithms for factorization of polynomials mod p which will give the individual irreducible polynomials (cf. the exercises) in practice. The importance of having irreducible polynomials at hand is that they give a representation of the finite fields \mathbb{F}_{p^n} (as quotients $\mathbb{F}_p[x]/(f(x))$ for $f(x)$ irreducible of degree n) conducive to explicit computations.

Note also that since the finite field \mathbb{F}_{p^n} is unique up to isomorphism, the quotients of $\mathbb{F}_p[x]$ by any of the irreducible polynomials of degree n are all isomorphic. If $f_1(x)$ and $f_2(x)$ are irreducible of degree n , then $f_2(x)$ splits completely in the field $\mathbb{F}_{p^n} \cong \mathbb{F}_p[x]/(f_1(x))$. If we denote a root of $f_2(x)$ by $\alpha(x)$ (to emphasize that it is a polynomial of degree $< n$ in x in $\mathbb{F}_p[x]/(f_1(x))$), then the isomorphism is given by

$$\begin{aligned}\mathbb{F}_p[x]/(f_2(x)) &\cong \mathbb{F}_p[x]/(f_1(x)) \\ x &\mapsto \alpha(x)\end{aligned}$$

(we have mapped a root of $f_2(x)$ in the first field to a root of $f_2(x)$ in the second field). For example, if $f_1(x) = x^4 + x^3 + 1$, $f_2(x) = x^4 + x + 1$ are two of the irreducible quartics over \mathbb{F}_2 determined above, then a simple computation verifies that

$$\alpha(x) = x^3 + x^2$$

is a root of $f_2(x)$ in $\mathbb{F}_{16} = \mathbb{F}_2[x]/(x^4 + x^3 + 1)$. Then we have

$$\begin{aligned}\mathbb{F}_2[x]/(x^4 + x + 1) &\cong \mathbb{F}_2[x]/(x^4 + x^3 + 1) \quad (\cong \mathbb{F}_{16}) \\ x &\mapsto x^3 + x^2.\end{aligned}$$

If we assume a result from elementary number theory we can give a formula for the number of irreducible polynomials of degree n . Define the *Möbius* μ -function by

$$\mu(n) = \begin{cases} 1 & \text{for } n = 1 \\ 0 & \text{if } n \text{ has a square factor} \\ (-1)^r & \text{if } n \text{ has } r \text{ distinct prime factors.} \end{cases}$$

If now $f(n)$ is a function defined for all nonnegative integers n and $F(n)$ is defined by

$$F(n) = \sum_{d|n} f(d) \quad n = 1, 2, \dots$$

then the *Möbius inversion formula* states that one can recover the function $f(n)$ from $F(n)$:

$$f(n) = \sum_{d|n} \mu(d) F\left(\frac{n}{d}\right) \quad n = 1, 2, \dots$$

This is an elementary result from number theory which we take for granted. Define

$$\psi(n) = \text{the number of irreducible polynomials of degree } n \text{ in } \mathbb{F}_p[x].$$

Counting degrees in Proposition 18 we have

$$p^n = \sum_{d|n} d\psi(d).$$

Applying the Möbius inversion formula (for $f(n) = n\psi(n)$) we obtain

$$n\psi(n) = \sum_{d|n} \mu(d) p^{n/d}$$

which gives us a formula for the number of irreducible polynomials of degree n over \mathbb{F}_p :

$$\psi(n) = \frac{1}{n} \sum_{d|n} \mu(d) p^{n/d}.$$

For example, in the case $p = 2, n = 4$ we have

$$\psi(4) = \frac{1}{4} [\mu(1)2^4 + \mu(2)2^2 + \mu(4)2^1] = \frac{1}{4}(16 - 4 + 0) = 3$$

as we determined directly above.

We have seen above that

$$\mathbb{F}_{p^m} \subseteq \mathbb{F}_{p^n} \text{ if and only if } m \text{ divides } n.$$

In particular, given any two finite fields $\mathbb{F}_{p^{n_1}}$ and $\mathbb{F}_{p^{n_2}}$ there is a third finite field containing (an isomorphic copy of) them, namely $\mathbb{F}_{p^{n_1 n_2}}$. This gives us a partial ordering on these fields and allows us to think of their union. Since these give *all* the finite extensions of \mathbb{F}_p , we see that the union of \mathbb{F}_{p^n} for all n is an algebraic closure of \mathbb{F}_p , unique up to isomorphism:

$$\overline{\mathbb{F}_p} = \bigcup_{n \geq 1} \mathbb{F}_{p^n}.$$

This provides a simple description of the algebraic closure of \mathbb{F}_p .

EXERCISES

1. Factor $x^8 - x$ into irreducibles in $\mathbb{Z}[x]$ and in $\mathbb{F}_2[x]$.
2. Write out the multiplication table for \mathbb{F}_4 and \mathbb{F}_8 .
3. Prove that an algebraically closed field must be infinite.
4. Construct the finite field of 16 elements and find a generator for the multiplicative group. How many generators are there?
5. Exhibit an explicit isomorphism between the splitting fields of $x^3 - x + 1$ and $x^3 - x - 1$ over \mathbb{F}_3 .
6. Suppose $K = \mathbb{Q}(\theta) = \mathbb{Q}(\sqrt{D_1}, \sqrt{D_2})$ with $D_1, D_2 \in \mathbb{Z}$, is a biquadratic extension and that $\theta = a + b\sqrt{D_1} + c\sqrt{D_2} + d\sqrt{D_1 D_2}$ where $a, b, c, d \in \mathbb{Z}$ are integers. Prove that the minimal polynomial $m_\theta(x)$ for θ over \mathbb{Q} is irreducible of degree 4 over \mathbb{Q} but is reducible modulo every prime p . In particular show that the polynomial $x^4 - 10x^2 + 1$ is irreducible in $\mathbb{Z}[x]$ but is reducible modulo every prime. [Use the fact that there are no biquadratic extensions over finite fields.]
7. Prove that one of 2, 3 or 6 is a square in \mathbb{F}_p for every prime p . Conclude that the polynomial
$$x^6 - 11x^4 + 36x^2 - 36 = (x^2 - 2)(x^2 - 3)(x^2 - 6)$$
has a root modulo p for every prime p but has no root in \mathbb{Z} .
8. Determine the splitting field of the polynomial $x^p - x - a$ over \mathbb{F}_p where $a \neq 0, a \in \mathbb{F}_p$. Show explicitly that the Galois group is cyclic. [Show $\alpha \mapsto \alpha + 1$ is an automorphism.] Such an extension is called an *Artin–Schreier extension* (cf. Exercise 9 of Section 7).
9. Let $q = p^m$ be a power of the prime p and let $\mathbb{F}_q = \mathbb{F}_{p^m}$ be the finite field with q elements. Let $\sigma_q = \sigma_p^m$ be the m^{th} power of the Frobenius automorphism σ_p , called the q -Frobenius automorphism.
 - Prove that σ_q fixes \mathbb{F}_q .
 - Prove that every finite extension of \mathbb{F}_q of degree n is the splitting field of $x^{q^n} - x$ over \mathbb{F}_q , hence is unique.
 - Prove that every finite extension of \mathbb{F}_q of degree n is cyclic with σ_q as generator.
 - Prove that the subfields of the unique extension of \mathbb{F}_q of degree n are in bijective correspondence with the divisors d of n .
10. Prove that n divides $\varphi(p^n - 1)$. [Observe that $\varphi(p^n - 1)$ is the order of the group of automorphisms of a cyclic group of order $p^n - 1$.]
11. Prove that $x^{p^n} - x + 1$ is irreducible over \mathbb{F}_p only when $n = 1$ or $n = p = 2$. [Note that if α is a root, then so is $\alpha + a$ for any $a \in \mathbb{F}_{p^n}$. Show that this implies $\mathbb{F}_p(\alpha)$ contains \mathbb{F}_{p^n} and that $[\mathbb{F}_p(\alpha) : \mathbb{F}_{p^n}] = p$.]

(*Berlekamp's Factorization Algorithm*) The following exercises outline the Berlekamp factorization algorithm for factoring polynomials in $\mathbb{F}_p[x]$. The efficiency of this algorithm is based on the efficiency of computing greatest common divisors in $\mathbb{F}_p[x]$ by the Euclidean Algorithm and on the efficiency of row-reduction matrix algorithms for solving systems of linear equations.

Let $f(x) \in \mathbb{F}_p[x]$ be a monic polynomial of degree n and let $f(x) = p_1(x)p_2(x)\dots p_k(x)$ where $p_1(x), p_2(x), \dots, p_k(x)$ are powers of distinct monic irreducibles in $\mathbb{F}_p[x]$.

12. Show that in order to write $f(x)$ as a product of irreducible polynomials in $\mathbb{F}_p[x]$ it suffices to determine the factors $p_1(x), \dots, p_k(x)$. [If $p(x) = q(x)^N \in \mathbb{F}_p[x]$ with $q(x)$ monic

and irreducible, show that $q(x)$ can be determined from $p(x)$ by checking for p^{th} powers and by computing greatest common divisors with derivatives.]

13. Let $g(x) \in \mathbb{F}_p[x]$ be any polynomial of degree $< n$. Denote by $R(h(x))$ the remainder of $h(x)$ after division by $f(x)$. Prove the following are equivalent:
 - (a) $R(g(x^p)) = g(x)$.
 - (b) $f(x)$ divides $[g(x)-0][g(x)-1]\dots[g(x)-(p-1)]$. [Use the fact that $g(x^p) = g(x)^p$ together with the factorization of $x^p - x$ in $\mathbb{F}_p[x]$.]
 - (c) $p_i(x)$ divides the product in (b) for $i = 1, 2, \dots, k$.
 - (d) For each i , $i = 1, 2, \dots, k$ there is an $s_i \in \mathbb{F}_p$ such that $p_i(x)$ divides $g(x) - s_i$, i.e., $g(x) \equiv s_i \pmod{p_i(x)}$.
14. Prove that the polynomials $g(x)$ of degree $< n$ satisfying the equivalent conditions of the previous exercise form a vector space V over \mathbb{F}_p of dimension k . [Use the Chinese Remainder Theorem applied to the p^k possible choices for the s_i in 13(d)].
15. Let $g(x) = b_0 + b_1x + \dots + b_{n-1}x^{n-1} \in V$. For $j = 0, 1, \dots, n-1$ let

$$R(x^{pj}) = a_{0,j} + a_{1,j}x + \dots + a_{n-1,j}x^{n-1}$$

and let A be the $n \times n$ matrix

$$A = \begin{pmatrix} a_{0,0} & a_{0,1} & \dots & a_{0,n-1} \\ a_{1,0} & a_{1,1} & \dots & a_{1,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n-1,0} & a_{n-1,1} & \dots & a_{n-1,n-1} \end{pmatrix}. \quad (*)$$

Show that condition (a) of Exercise 13 for $g(x) \in V$ is equivalent to

$$(A - I)B = 0 \quad (**)$$

where B is the column matrix with entries b_0, b_1, \dots, b_{n-1} . Conclude that the rank of the matrix $A - I$ is $n - k$. Note that this already suffices to determine if $f(x)$ is irreducible, without actually determining the factors.

16. Let $g_1(x), g_2(x), \dots, g_k(x)$ be a basis of solutions to $(**)$ (so a basis for V), where we may take $g_1(x) = 1$. Beginning with $w(x) = f(x)$, compute the greatest common divisor $(w(x), g_i(x) - s)$ for $i = 2, 3, \dots, k$ and $s \in \mathbb{F}_p$ for every factor of $f(x)$ already computed. Note by Exercise 13(d) that every factor $p_i(x)$ of $f(x)$ divides such a g.c.d. The process terminates when k relatively prime factors have been determined.

Prove that this procedure actually gives all the factors $p_1(x), p_2(x), \dots, p_k(x)$, i.e., one can separate the individual factors $p_1(x), p_2(x), \dots, p_k(x)$ by this procedure, as follows:

If this were not the case, then for two of the factors, say $p_1(x)$ and $p_2(x)$, for each $i = 1, 2, \dots, k$ there would exist $s_i \in \mathbb{F}_p$ such that $g_i(x) - s_i$ is divisible by both $p_1(x)$ and $p_2(x)$. By the Chinese Remainder Theorem, choose a $g(x) \in V$ satisfying $g(x) \equiv 0 \pmod{p_1(x)}$ and $g(x) \equiv 1 \pmod{p_2(x)}$. Write $g(x) = \sum_{i=1}^k c_i g_i(x)$ in terms of the basis for V and let $s = \sum_{i=1}^k c_i s_i(x) \in \mathbb{F}_p$. Show that $s \equiv 0 \pmod{p_1(x)}$ so that $s = 0$ and $s \equiv 1 \pmod{p_2(x)}$ so that $s = 1$, a contradiction.

17. This exercise follows Berlekamp's Factorization Algorithm outlined in the previous exercises to determine the factorization of $f(x) = x^5 + x^2 + 4x + 6$ in $\mathbb{F}_7[x]$.
 - (a) Show that $x^7 \equiv x^2 + 3x^3 + 6x^4 \pmod{f(x)}$. Similarly compute x^{14}, x^{21} , and x^{28} modulo $f(x)$ (note that x^{14} can most easily be computed by squaring the result for

x^7 and then reducing, etc.) to show that in this case the matrix A in Exercise 15 is

$$\begin{pmatrix} 1 & 0 & 5 & 1 & 4 \\ 0 & 0 & 1 & 1 & 2 \\ 0 & 1 & 3 & 3 & 3 \\ 0 & 3 & 4 & 2 & 2 \\ 0 & 6 & 3 & 1 & 1 \end{pmatrix}.$$

- (b) Show that the reduced row echelon form for $A - I$ is the matrix

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 6 \\ 0 & 0 & 1 & 0 & 6 \\ 0 & 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Conclude that $k = 2$ (so $f(x)$ is the product of precisely two factors which are powers of irreducible polynomials) and that $g_1(x) = 1$ and $g_2(x) = x^4 + 5x^3 + x^2 + x$ give a basis for the solutions to $(**)$ in Exercise 15.

- (c) Following the procedure in Exercise 16, show that $(f(x), g_2(x) - 1) = x^2 + 3x + 5 = p_1(x)$, with $f(x)/p_1(x) = x^3 + 4x^2 + 4x + 4 = p_2(x)$, giving the powers of the irreducible polynomials dividing $f(x)$ in $\mathbb{F}_7[x]$. Show that neither factor is a 7th power in $\mathbb{F}_7[x]$ and that each is relatively prime to its derivative to conclude that both factors are irreducible polynomials, giving the complete factorization of $f(x)$ into irreducible polynomials:

$$f(x) = (x^2 + 3x + 5)(x^3 + 4x^2 + 4x + 4) \in \mathbb{F}_7[x].$$

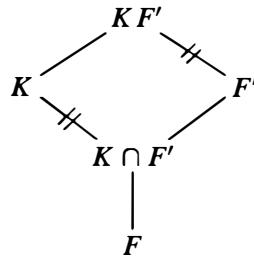
14.4 COMPOSITE EXTENSIONS AND SIMPLE EXTENSIONS

We now consider the effect of taking composites with Galois extensions. The first result states that “sliding up” a Galois extension gives a Galois extension.

Proposition 19. Suppose K/F is a Galois extension and F'/F is any extension. Then KF'/F' is a Galois extension, with Galois group

$$\text{Gal}(KF'/F') \cong \text{Gal}(K/K \cap F')$$

isomorphic to a subgroup of $\text{Gal}(K/F)$. Pictorially,



Proof: If K/F is Galois, then K is the splitting field of some separable polynomial $f(x)$ in $F[x]$. Then KF'/F' is the splitting field of $f(x)$ viewed as a polynomial in

$F'[x]$, hence this extension is Galois. Since K/F is Galois, every embedding of K fixing F is an automorphism of K , so the map

$$\begin{aligned}\varphi : \text{Gal}(KF'/F') &\rightarrow \text{Gal}(K/F) \\ \sigma &\mapsto \sigma|_K\end{aligned}$$

defined by restricting an automorphism σ to the subfield K is well defined. It is clearly a homomorphism, with kernel

$$\ker \varphi = \{\sigma \in \text{Gal}(KF'/F') \mid \sigma|_K = 1\}.$$

Since an element in $\text{Gal}(KF'/F')$ is trivial on F' , the elements in the kernel are trivial both on K and on F' , hence on their composite, so the kernel consists only of the identity automorphism. Hence φ is injective.

Let H denote the image of φ in $\text{Gal}(K/F)$ and let K_H denote the corresponding fixed subfield of K containing F . Since every element in H fixes F' , K_H contains $K \cap F'$. On the other hand, the composite $K_H F'$ is fixed by $\text{Gal}(KF'/F')$ (any $\sigma \in \text{Gal}(KF'/F')$ fixes F' and acts on $K_H \subseteq K$ via its restriction $\sigma|_K \in H$, which fixes K_H by definition). By the Fundamental Theorem it follows that $K_H F' = F'$, so that $K_H \subseteq F'$, which gives the reverse inclusion $K_H \subseteq K \cap F'$. Hence $K_H = K \cap F'$, so again by the Fundamental Theorem, $H = \text{Gal}(K/K \cap F')$, completing the proof.

Corollary 20. Suppose K/F is a Galois extension and F'/F is any finite extension. Then

$$[KF' : F] = \frac{[K : F][F' : F]}{[K \cap F' : F]}.$$

Proof: This follows by the proposition from the equality $[KF' : F'] = [K : K \cap F']$ given by the orders of the Galois groups in the proposition.

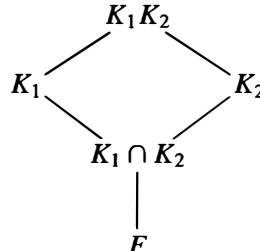
The example $F = \mathbb{Q}$, $K = \mathbb{Q}(\sqrt[3]{2})$, $F' = \mathbb{Q}(\rho\sqrt[3]{2})$, ρ a primitive 3rd root of unity, shows that the formula of Corollary 20 does not hold in general if neither of the two extensions is Galois.

Proposition 21. Let K_1 and K_2 be Galois extensions of a field F . Then

- (1) The intersection $K_1 \cap K_2$ is Galois over F .
- (2) The composite $K_1 K_2$ is Galois over F . The Galois group is isomorphic to the subgroup

$$H = \{(\sigma, \tau) \mid \sigma|_{K_1 \cap K_2} = \tau|_{K_1 \cap K_2}\}$$

of the direct product $\text{Gal}(K_1/F) \times \text{Gal}(K_2/F)$ consisting of elements whose restrictions to the intersection $K_1 \cap K_2$ are equal.



Proof: (1) Suppose $p(x)$ is an irreducible polynomial in $F[x]$ with a root α in $K_1 \cap K_2$. Since $\alpha \in K_1$ and K_1/F is Galois, all the roots of $p(x)$ lie in K_1 . Similarly all the roots lie in K_2 , hence all the roots of $p(x)$ lie in $K_1 \cap K_2$. It follows easily that $K_1 \cap K_2$ is Galois as in Theorem 13.

(2) If K_1 is the splitting field of the separable polynomial $f_1(x)$ and K_2 is the splitting field of the separable polynomial $f_2(x)$ then the composite is the splitting field for the squarefree part of the polynomial $f_1(x)f_2(x)$, hence is Galois over F .

The map

$$\begin{aligned}\varphi : \text{Gal}(K_1 K_2 / F) &\rightarrow \text{Gal}(K_1 / F) \times \text{Gal}(K_2 / F) \\ \sigma &\mapsto (\sigma|_{K_1}, \sigma|_{K_2})\end{aligned}$$

is clearly a homomorphism. The kernel consists of the elements σ which are trivial on both K_1 and K_2 , hence trivial on the composite, so the map is injective. The image lies in the subgroup H , since

$$(\sigma|_{K_1})|_{K_1 \cap K_2} = \sigma|_{K_1 \cap K_2} = (\sigma|_{K_2})|_{K_1 \cap K_2}.$$

The order of H can be computed by observing that for every $\sigma \in \text{Gal}(K_1 / F)$ there are $|\text{Gal}(K_2 / K_1 \cap K_2)|$ elements $\tau \in \text{Gal}(K_2 / F)$ whose restrictions to $K_1 \cap K_2$ are $\sigma|_{K_1 \cap K_2}$. Hence

$$\begin{aligned}|H| &= |\text{Gal}(K_1 / F)| \cdot |\text{Gal}(K_2 / K_1 \cap K_2)| \\ &= |\text{Gal}(K_1 / F)| \frac{|\text{Gal}(K_2 / F)|}{|\text{Gal}(K_1 \cap K_2 / F)|}.\end{aligned}$$

By Corollary 20 and the diagram above we see that the orders of H and $\text{Gal}(K_1 K_2 / F)$ are then both equal to

$$[K_1 K_2 : F] = \frac{[K_1 : F][K_2 : F]}{[K_1 \cap K_2 : F]}.$$

Hence the image of φ is precisely H , completing the proof.

Corollary 22. Let K_1 and K_2 be Galois extensions of a field F with $K_1 \cap K_2 = F$. Then

$$\text{Gal}(K_1 K_2 / F) \cong \text{Gal}(K_1 / F) \times \text{Gal}(K_2 / F).$$

Conversely, if K is Galois over F and $G = \text{Gal}(K / F) = G_1 \times G_2$ is the direct product of two subgroups G_1 and G_2 , then K is the composite of two Galois extensions K_1 and K_2 of F with $K_1 \cap K_2 = F$.

Proof: The first part follows immediately from the proposition. For the second, let K_1 be the fixed field of $G_1 \subset G$ and let K_2 be the fixed field of $G_2 \subset G$. Then $K_1 \cap K_2$ is the field corresponding to the subgroup $G_1 G_2$, which is all of G in this case, so $K_1 \cap K_2 = F$. The composite $K_1 K_2$ is the field corresponding to the subgroup $G_1 \cap G_2$, which is the identity here, so $K_1 K_2 = K$, completing the proof.

Corollary 23. Let E/F be any finite separable extension. Then E is contained in an extension K which is Galois over F and is minimal in the sense that in a fixed algebraic closure of K any other Galois extension of F containing E contains K .

Proof: There exists a Galois extension of F containing E , for example the composite of the splitting fields of the minimal polynomials for a basis for E over F (which are all separable since E is separable over F). Then the intersection of all the Galois extensions of F containing E is the field K .

Definition. The Galois extension K of F containing E in the previous corollary is called the *Galois closure* of E over F .

It is often simpler to work in a Galois extension (for example in computing degrees as in Corollary 20). The existence of a Galois closure for a separable extension is frequently useful for reducing computations to consideration of Galois extensions.

Recall that an extension K of F is called *simple* if $K = F(\theta)$ for some element θ , in which case θ is called a *primitive element* for K .

Proposition 24. Let K/F be a finite extension. Then $K = F(\theta)$ if and only if there exist only finitely many subfields of K containing F .

Proof: Suppose first that $K = F(\theta)$ is simple. Let E be a subfield of K containing F : $F \subseteq E \subseteq K$. Let $f(x) \in F[x]$ be the minimal polynomial for θ over F and let $g(x) \in E[x]$ be the minimal polynomial for θ over E . Then $g(x)$ divides $f(x)$ in $E[x]$. Let E' be the field generated over F by the coefficients of $g(x)$. Then $E' \subseteq E$ and clearly the minimal polynomial for θ over E' is still $g(x)$. But then

$$[K : E] = \deg g(x) = [K : E']$$

implies that $E = E'$. It follows that the subfields of K containing F are the subfields generated by the coefficients of the monic factors of $f(x)$, hence there are finitely many such subfields.

Suppose conversely that there are finitely many subfields of K containing F . If F is a finite field, then we have already seen that K is a simple extension (Proposition 17). Hence we may suppose F is infinite. It clearly suffices to show that $F(\alpha, \beta)$ is generated by a single element since K is finitely generated over F . Consider the subfields

$$F(\alpha + c\beta), \quad c \in F.$$

Then since there are infinitely many choices for $c \in F$ and only finitely many such subfields, there exist c, c' in F , $c \neq c'$, with

$$F(\alpha + c\beta) = F(\alpha + c'\beta).$$

Then $\alpha + c\beta$ and $\alpha + c'\beta$ both lie in $F(\alpha + c\beta)$, and taking their difference shows that $(c - c')\beta \in F(\alpha + c\beta)$. Hence $\beta \in F(\alpha + c\beta)$ and then also $\alpha \in F(\alpha + c\beta)$. Therefore $F(\alpha, \beta) \subseteq F(\alpha + c\beta)$ and since the reverse inclusion is obvious, we have

$$F(\alpha, \beta) = F(\alpha + c\beta),$$

completing the proof.

Theorem 25. (The Primitive Element Theorem) If K/F is finite and separable, then K/F is simple. In particular, any finite extension of fields of characteristic 0 is simple.

Proof: Let L be the Galois closure of K over F . Then any subfield of K containing F corresponds to a subgroup of the Galois group $\text{Gal}(L/F)$ by the Fundamental Theorem. Since there are only finitely many such subgroups, the previous proposition shows that K/F is simple. The last statement follows since any finite extension of fields in characteristic 0 is separable.

As the proof of the proposition indicates, a primitive element for an extension can be obtained as a simple linear combination of the generators for the extension. In the case of Galois extensions it is only necessary to determine a linear combination which is not fixed by any nontrivial element of the Galois group since then by the Fundamental Theorem this linear combination could not lie in any proper subfield.

Examples

- (1) The element $\sqrt{2} + \sqrt{3}$ generates the field $\mathbb{Q}(\sqrt{2}, \sqrt{3})$ as we have already seen (it is not fixed by any of the four Galois automorphisms of this field).
- (2) The field $\overline{\mathbb{F}_p}(x, y)$ of rational functions in the variables x and y over the algebraic closure $\overline{\mathbb{F}_p}$ of \mathbb{F}_p is not a simple extension of the subfield $F = \overline{\mathbb{F}_p}(x^p, y^p)$. It is easy to see that

$$[\overline{\mathbb{F}_p}(x, y) : \overline{\mathbb{F}_p}(x^p, y^p)] = p^2$$

and that the subfields

$$F(x + cy), \quad c \in \overline{\mathbb{F}_p}$$

are all of degree p over $\overline{\mathbb{F}_p}(x^p, y^p)$ (note that $(x + cy)^p = x^p + c^p y^p \in \overline{\mathbb{F}_p}(x^p, y^p)$). If any two of these subfields were equal, then just as in the proof of Proposition 24 we would have

$$\overline{\mathbb{F}_p}(x, y) = F(x + cy)$$

which is impossible by degree considerations. Hence there are infinitely many such subfields and the extension cannot be simple.

EXERCISES

1. Determine the Galois closure of the field $\mathbb{Q}(\sqrt{1 + \sqrt{2}})$ over \mathbb{Q} .
2. Find a primitive generator for $\mathbb{Q}(\sqrt{2}, \sqrt{3}, \sqrt{5})$ over \mathbb{Q} .
3. Let F be a field contained in the ring of $n \times n$ matrices over \mathbb{Q} . Prove that $[F : \mathbb{Q}] \leq n$. (Note that, by Exercise 19 of Section 13.2, the ring of $n \times n$ matrices over \mathbb{Q} does contain fields of degree n over \mathbb{Q} .)
4. Let $f(x) \in F[x]$ be an irreducible polynomial of degree n over the field F , let L be the splitting field of $f(x)$ over F and let α be a root of $f(x)$ in L . If K is any Galois extension of F , show that the polynomial $f(x)$ splits into a product of m irreducible polynomials each of degree d over K , where $d = [K(\alpha) : K] = [(L \cap K)(\alpha) : L \cap K]$ and $m = n/d = [F(\alpha) \cap K : F]$. [Show first that the factorization of $f(x)$ over K is the same as its factorization over $L \cap K$. Then if H is the subgroup of the Galois group of L

over F corresponding to $L \cap K$ the factors of $f(x)$ over $L \cap K$ correspond to the orbits of H on the roots of $f(x)$. Use Exercise 9 of Section 4.1.]

5. Let p be a prime and let F be a field. Let K be a Galois extension of F whose Galois group is a p -group (i.e., the degree $[K : F]$ is a power of p). Such an extension is called a *p -extension* (note that p -extensions are Galois by definition).
 - (a) Let L be a p -extension of K . Prove that the Galois closure of L over F is a p -extension of F .
 - (b) Give an example to show that (a) need not hold if $[K : F]$ is a power of p but K/F is not Galois.
6. Prove that $\mathbb{F}_p(x, y)/\mathbb{F}_p(x^p, y^p)$ is not a simple extension by explicitly exhibiting an infinite number of intermediate subfields.
7. Let $F \subseteq K \subseteq L$ and let $\theta \in L$ with $p(x) = m_{\theta, F}(x)$. Prove that $K \otimes_F F(\theta) \cong K[x]/(p(x))$ as K -algebras.
8. Let K_1 and K_2 be two algebraic extensions of a field F contained in the field L of characteristic zero. Prove that the F -algebra $K_1 \otimes_F K_2$ has no nonzero nilpotent elements. [Use the preceding exercise.]

14.5 CYCLOTOMIC EXTENSIONS AND ABELIAN EXTENSIONS OVER \mathbb{Q}

We have already determined that the cyclotomic field $\mathbb{Q}(\zeta_n)$ of n^{th} roots of unity is a Galois extension of \mathbb{Q} of degree $\varphi(n)$ where φ denotes the Euler φ -function. Any automorphism of this field is uniquely determined by its action on the primitive n^{th} root of unity ζ_n . This element must be mapped to another primitive n^{th} root of unity (recall these are the roots of the irreducible cyclotomic polynomial $\Phi_n(x)$). Hence $\sigma(\zeta_n) = \zeta_n^a$ for some integer a , $1 \leq a < n$, relatively prime to n . Since there are precisely $\varphi(n)$ such integers a it follows that in fact each of these maps is indeed an automorphism of $\mathbb{Q}(\zeta_n)$. Note also that we can define σ_a for any integer a relatively prime to n by the same formula and that σ_a depends only on the residue class of a modulo n .

Theorem 26. The Galois group of the cyclotomic field $\mathbb{Q}(\zeta_n)$ of n^{th} roots of unity is isomorphic to the multiplicative group $(\mathbb{Z}/n\mathbb{Z})^\times$. The isomorphism is given explicitly by the map

$$\begin{aligned} (\mathbb{Z}/n\mathbb{Z})^\times &\xrightarrow{\sim} \text{Gal}(\mathbb{Q}(\zeta_n)/\mathbb{Q}) \\ a \pmod{n} &\mapsto \sigma_a \end{aligned}$$

where σ_a is the automorphism defined by

$$\sigma_a(\zeta_n) = \zeta_n^a.$$

Proof: The discussion above shows that σ_a is an automorphism for any $a \pmod{n}$, so the map above is well defined. It is a homomorphism since

$$\begin{aligned} (\sigma_a \sigma_b)(\zeta_n) &= \sigma_a(\zeta_n^b) = (\zeta_n^b)^a \\ &= \zeta_n^{ab} \end{aligned}$$

which shows that $\sigma_a \sigma_b = \sigma_{ab}$. The map is bijective by the discussion above since we know that every Galois automorphism is of the form σ_a for a uniquely defined $a \pmod{n}$. Hence the map is an isomorphism.

Examples

- (1) The field $\mathbb{Q}(\zeta_5)$ is Galois over \mathbb{Q} with Galois group $(\mathbb{Z}/5\mathbb{Z})^\times \cong \mathbb{Z}/4\mathbb{Z}$. This is our first example of a Galois extension of \mathbb{Q} of degree 4 with a *cyclic* Galois group. The elements of the Galois group are $\{\sigma_1 = 1, \sigma_2, \sigma_3, \sigma_4\}$ in the notation above. A generator for this cyclic group is $\sigma_2 : \zeta_5 \mapsto \zeta_5^2$ (since 2 has order 4 in $(\mathbb{Z}/5\mathbb{Z})^\times$).

There is precisely one nontrivial subfield, a quadratic extension of \mathbb{Q} , the fixed field of the subgroup $\{1, \sigma_4 = \sigma_{-1}\}$. An element in this subfield is given by

$$\alpha = \zeta_5 + \sigma_{-1}\zeta_5 = \zeta_5 + \zeta_5^{-1}$$

since this element is clearly fixed by σ_{-1} . The element ζ_5 satisfies

$$\zeta_5^4 + \zeta_5^3 + \zeta_5^2 + \zeta_5 + 1 = 0.$$

Notice then that

$$\begin{aligned} \alpha^2 + \alpha - 1 &= (\zeta_5^2 + 2 + \zeta_5^{-2}) + (\zeta_5 + \zeta_5^{-1}) - 1 \\ &= \zeta_5^2 + 2 + \zeta_5^3 + \zeta_5 + \zeta_5^4 - 1 = 0. \end{aligned}$$

Solving explicitly for α we see that the quadratic extension of \mathbb{Q} generated by α is $\mathbb{Q}(\sqrt{5})$:

$$\mathbb{Q}(\zeta_5 + \zeta_5^{-1}) = \mathbb{Q}(\sqrt{5}).$$

It can be shown in general (this is not completely trivial) that for p an odd prime the field $\mathbb{Q}(\zeta_p)$ contains the quadratic field $\mathbb{Q}(\sqrt{\pm p})$, where the $+$ sign is correct if $p \equiv 1 \pmod{4}$ and the $-$ sign is correct if $p \equiv 3 \pmod{4}$ (cf. Exercise 11 in Section 7).

- (2) $\mathbb{Q}(\zeta_{13})$. For p an odd prime we can construct a primitive element for any of the subfields of $\mathbb{Q}(\zeta_p)$ as in the previous example. A basis for $\mathbb{Q}(\zeta_p)$ over \mathbb{Q} is given by

$$1, \zeta_p, \zeta_p^2, \dots, \zeta_p^{p-2}.$$

Since

$$\zeta_p^{p-1} + \zeta_p^{p-2} + \dots + \zeta_p + 1 = 0$$

we see that also the elements

$$\zeta_p, \zeta_p^2, \dots, \zeta_p^{p-2}, \zeta_p^{p-1}$$

form a basis. The reason for choosing this basis is that any σ in the Galois group $\text{Gal}(\mathbb{Q}(\zeta_p)/\mathbb{Q})$ simply *permutes* these basis elements since these are precisely the primitive p^{th} roots of unity. Note that it is at this point that we need p to be a prime — in general the primitive n^{th} roots of unity do not give a basis for the cyclotomic field of n^{th} roots of unity over \mathbb{Q} (for example, the primitive 4^{th} roots of unity, $\pm i$, are not linearly independent).

Let H be any subgroup of the Galois group of $\mathbb{Q}(\zeta_p)$ over \mathbb{Q} and let

$$\alpha_H = \sum_{\sigma \in H} \sigma \zeta_p, \tag{14.10}$$

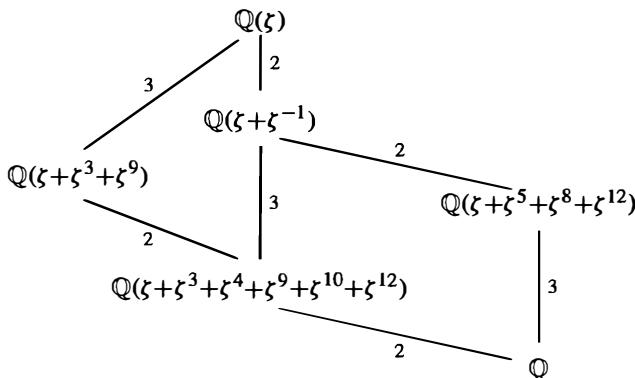
the sum of the conjugates of ζ_p by the elements in H . For any $\tau \in H$, the elements $\tau\alpha$ run over the elements of H as σ runs over the elements of H . It follows that $\tau\alpha = \alpha$, so

that α lies in the fixed field for H . If now τ is *not* an element of H , then $\tau\alpha$ is the sum of basis elements (recall that any automorphism permutes the basis elements here), one of which is $\tau(\zeta_p)$. If we had $\tau\alpha = \alpha$ then since these elements are a basis, we must have $\tau(\zeta_p) = \sigma(\zeta_p)$ for one of the terms $\sigma\zeta_p$ in (10). But this implies $\tau\sigma^{-1} = 1$ since this automorphism is the identity on ζ_p . Then $\tau = \sigma \in H$, a contradiction. This shows that α is not fixed by any automorphism not contained in H , so that $\mathbb{Q}(\alpha)$ is precisely the fixed field of H .

For a specific example, consider the subfields of $\mathbb{Q}(\zeta_{13})$, which correspond to the subgroups of $(\mathbb{Z}/13\mathbb{Z})^\times \cong \mathbb{Z}/12\mathbb{Z}$. A generator for this cyclic group is the automorphism $\sigma = \sigma_2$ which maps ζ_{13} to ζ_{13}^2 . The nontrivial subgroups correspond to the nontrivial divisors of 12, hence are of orders 2, 3, 4, and 6 with generators $\sigma^6, \sigma^4, \sigma^3$ and σ^2 , respectively. The corresponding fixed fields will be of degrees 6, 4, 3 and 2 over \mathbb{Q} , respectively. Generators are given by ($\zeta = \zeta_{13}$)

$$\begin{aligned}\zeta + \sigma^6\zeta &= \zeta + \zeta^{2^6} = \zeta + \zeta^{-1} \\ \zeta + \sigma^4\zeta + \sigma^8\zeta &= \zeta + \zeta^{2^4} + \zeta^{2^8} = \zeta + \zeta^3 + \zeta^9 \\ \zeta + \sigma^3\zeta + \sigma^6\zeta + \sigma^9\zeta &= \zeta + \zeta^8 + \zeta^{12} + \zeta^5 \\ \zeta + \sigma^2\zeta + \sigma^4\zeta + \sigma^6\zeta + \sigma^8\zeta + \sigma^{10}\zeta &= \zeta + \zeta^4 + \zeta^3 + \zeta^{12} + \zeta^9 + \zeta^{10}.\end{aligned}$$

The lattice of subfields for this extension is the following:



The elements constructed in equation (10) and their conjugates are called the *periods* of ζ and are useful in the study of the arithmetic of the cyclotomic fields. The study of their combinatorial properties is referred to as *cyclotomy*.

Suppose that $n = p_1^{a_1} p_2^{a_2} \cdots p_k^{a_k}$ is the decomposition of n into distinct prime powers. Since $\zeta_n^{p_2^{a_2} \cdots p_k^{a_k}}$ is a primitive $p_1^{a_1}$ -th root of unity, the field $K_1 = \mathbb{Q}(\zeta_{p_1^{a_1}})$ is a subfield of $\mathbb{Q}(\zeta_n)$. Similarly, each of the fields $K_i = \mathbb{Q}(\zeta_{p_i^{a_i}})$, $i = 1, 2, \dots, k$ is a subfield of $\mathbb{Q}(\zeta_n)$. The composite of the fields contains the product $\zeta_{p_1^{a_1}} \zeta_{p_2^{a_2}} \cdots \zeta_{p_k^{a_k}}$, which is a primitive n^{th} root of unity, hence the composite field is $\mathbb{Q}(\zeta_n)$. Since the extension degrees $[K_i : \mathbb{Q}]$ equal $\varphi(p_i^{a_i})$, $i = 1, 2, \dots, k$ and $\varphi(n) = \varphi(p_1^{a_1})\varphi(p_2^{a_2}) \cdots \varphi(p_k^{a_k})$, the degree of the composite of the fields K_i is precisely the product of the degrees of the K_i . It follows from Proposition 21 (and a simple induction from the two fields considered in the proposition to the k fields here) that the intersection of all these fields

is precisely \mathbb{Q} . Then Corollary 22 shows that the Galois group for $\mathbb{Q}(\zeta_n)$ is the direct product of the Galois groups over \mathbb{Q} for the subfields K_i . We summarize this as the following corollary.

Corollary 27. Let $n = p_1^{a_1} p_2^{a_2} \cdots p_k^{a_k}$ be the decomposition of the positive integer n into distinct prime powers. Then the cyclotomic fields $\mathbb{Q}(\zeta_{p_i^{a_i}})$, $i = 1, 2, \dots, k$ intersect only in the field \mathbb{Q} and their composite is the cyclotomic field $\mathbb{Q}(\zeta_n)$. We have

$$\text{Gal}(\mathbb{Q}(\zeta_n)/\mathbb{Q}) \cong \text{Gal}(\mathbb{Q}(\zeta_{p_1^{a_1}})/\mathbb{Q}) \times \text{Gal}(\mathbb{Q}(\zeta_{p_2^{a_2}})/\mathbb{Q}) \times \cdots \times \text{Gal}(\mathbb{Q}(\zeta_{p_k^{a_k}})/\mathbb{Q})$$

which under the isomorphism in Theorem 26 is the Chinese Remainder Theorem:

$$(\mathbb{Z}/n\mathbb{Z})^\times \cong (\mathbb{Z}/p_1^{a_1}\mathbb{Z})^\times \times (\mathbb{Z}/p_2^{a_2}\mathbb{Z})^\times \times \cdots \times (\mathbb{Z}/p_k^{a_k}\mathbb{Z})^\times.$$

Proof: The only statement which has not been proved is the identification of the isomorphism of Galois groups with the statement of the Chinese Remainder Theorem on the group $(\mathbb{Z}/n\mathbb{Z})^\times$, which is quite simple and is left for the exercises.

By Theorem 26 the Galois group of $\mathbb{Q}(\zeta_n)/\mathbb{Q}$ is in particular an abelian group.

Definition. The extension K/F is called an *abelian* extension if K/F is Galois and $\text{Gal}(K/F)$ is an abelian group.

Since all the subgroups and quotient groups of abelian groups are abelian, we see by the Fundamental Theorem of Galois Theory that every subfield containing F of an abelian extension of F is again an abelian extension of F . By the results on composites of extensions in the last section, we also see that the composite of abelian extensions is again an abelian extension (since the Galois group of the composite is isomorphic to a subgroup of the direct product of the Galois groups, hence is abelian).

It is an open problem to determine which groups arise as the Galois groups of Galois extensions of \mathbb{Q} . Using the results above we can see that every *abelian* group appears as the Galois group of some extension of \mathbb{Q} , in fact as the Galois group of some subfield of a cyclotomic field.

Let $n = p_1 p_2 \cdots p_k$ be the product of distinct primes. Then by the Chinese Remainder Theorem

$$\begin{aligned} (\mathbb{Z}/n\mathbb{Z})^\times &\cong (\mathbb{Z}/p_1\mathbb{Z})^\times \times (\mathbb{Z}/p_2\mathbb{Z})^\times \times \cdots \times (\mathbb{Z}/p_k\mathbb{Z})^\times \\ &\cong Z_{p_1-1} \times Z_{p_2-1} \times \cdots \times Z_{p_k-1}. \end{aligned} \tag{14.11}$$

Now, suppose G is any finite abelian group. By the Fundamental Theorem for Abelian Groups,

$$G \cong Z_{n_1} \times Z_{n_2} \times \cdots \times Z_{n_k}$$

for some integers n_1, n_2, \dots, n_k . We take as known that given any integer m there are infinitely many primes p with $p \equiv 1 \pmod{m}$ (see the exercises following Section 13.6

for one proof using cyclotomic polynomials). Given this result, choose distinct primes p_1, p_2, \dots, p_k such that

$$p_1 \equiv 1 \pmod{n_1}$$

$$p_2 \equiv 1 \pmod{n_2}$$

$$\vdots$$

$$p_k \equiv 1 \pmod{n_k}$$

and let $n = p_1 p_2 \cdots p_k$ as above.

By construction, n_i divides $p_i - 1$ for $i = 1, 2, \dots, k$, so the group \mathbb{Z}_{p_i-1} has a subgroup H_i of order $\frac{p_i - 1}{n_i}$ for $i = 1, 2, \dots, k$, and the quotient by this subgroup is cyclic of order n_i . Hence the quotient of $(\mathbb{Z}/n\mathbb{Z})^\times$ in equation (11) by $H_1 \times H_2 \times \cdots \times H_k$ is isomorphic to the group G .

By Theorem 26 and the Fundamental Theorem of Galois Theory, we see that there is a subfield of $\mathbb{Q}(\zeta_{p_1 p_2 \cdots p_k})$ which is Galois over \mathbb{Q} with G as Galois group. We summarize this in the following corollary.

Corollary 28. Let G be any finite abelian group. Then there is a subfield K of a cyclotomic field with $\text{Gal}(K/\mathbb{Q}) \cong G$.

There is a converse to this result (whose proof is beyond our scope), the celebrated Kronecker–Weber Theorem:

Theorem (Kronecker–Weber) Let K be a finite abelian extension of \mathbb{Q} . Then K is contained in a cyclotomic extension of \mathbb{Q} .

The abelian extensions of \mathbb{Q} are the “easiest” Galois extensions (at least in so far as the structure of their Galois groups is concerned) and the previous result shows they can be classified by the cyclotomic extensions of \mathbb{Q} . For other finite extensions of \mathbb{Q} as base field, it is more difficult to describe the abelian extensions. The study of the abelian extensions of an arbitrary finite extension F of \mathbb{Q} is referred to as *class field theory*. There is a classification of the abelian extensions of F by invariants associated to F which greatly generalizes the results on cyclotomic fields over \mathbb{Q} . In general, however, the construction of abelian extensions is not nearly as explicit as in the case of the cyclotomic fields. One case where such a description is possible is for the abelian extensions of an imaginary quadratic field $(\mathbb{Q}(\sqrt{-D}))$ for D positive), where the abelian extensions can be constructed by adjoining values of certain elliptic functions (this is the analogue of adjoining the roots of unity, which are the values of the exponential function e^x for certain x). The study of the arithmetic of such abelian extensions and the search for similar results for non-abelian extensions are rich and fascinating areas of current mathematical research.

We end our discussion of the cyclotomic fields with the problem of the constructibility of the regular n -gon by straightedge and compass.

Recall (cf. Section 13.3) that an element α is constructible over \mathbb{Q} if and only if the field $\mathbb{Q}(\alpha)$ is contained in a field K obtained by a series of quadratic extensions:

$$\mathbb{Q} = K_0 \subset K_1 \subset \cdots \subset K_i \subset K_{i+1} \subset \cdots \subset K_m = K \quad (14.12)$$

with

$$[K_{i+1} : K_i] = 2, \quad i = 0, 1, \dots, m-1.$$

The construction of the regular n -gon in \mathbb{R}^2 is evidently equivalent to the construction of the n^{th} roots of unity, since the n^{th} roots of unity form the vertices of a regular n -gon on the unit circle in \mathbb{C} with one vertex at the point 1.

The construction of ζ_n is equivalent to the constructibility of the first coordinate x in \mathbb{R}^2 of ζ_n , namely the real part of ζ_n . Since the complex conjugate of ζ_n is just ζ_n^{-1} , the real part of ζ_n is $x = \frac{1}{2}(\zeta_n + \zeta_n^{-1})$. Note that ζ_n satisfies the quadratic equation $\zeta_n^2 - 2x\zeta_n + 1 = 0$ over $\mathbb{Q}(x)$. Since $\mathbb{Q}(x)$ consists only of real numbers, it follows that $[\mathbb{Q}(\zeta_n) : \mathbb{Q}(x)] = 2$, so that $\mathbb{Q}(x)$ is an extension of degree $\varphi(n)/2$ of \mathbb{Q} .

It follows that if the regular n -gon can be constructed by straightedge and compass then $\varphi(n)$ must be a power of 2. Conversely, if $\varphi(n) = 2^m$ is a power of 2, then the Galois group $\text{Gal}(\mathbb{Q}(\zeta_n)/\mathbb{Q})$ is an abelian group whose order is a power of 2, so the same is true for the Galois group $\text{Gal}(\mathbb{Q}(x)/\mathbb{Q})$. It is easy to see by the Fundamental Theorem for Abelian Groups that an abelian group G of order 2^m has a chain of subgroups

$$G = G_m > G_{m-1} > \cdots > G_{i+1} > G_i > \cdots > G_0 = 1$$

with

$$[G_{i+1} : G_i] = 2, \quad i = 0, 1, 2, \dots, m-1.$$

Applying this to the group $G = \text{Gal}(\mathbb{Q}(x)/\mathbb{Q})$ and taking the fixed fields for the subgroups G_i , $i = 0, 1, \dots, m-1$, we obtain (by the Fundamental Theorem of Galois Theory) a sequence of quadratic extensions as in (12) above.

We conclude that the regular n -gon can be constructed by straightedge and compass if and only if $\varphi(n)$ is a power of 2. Decomposing n into prime powers to compute $\varphi(n)$ we see that this means $n = 2^k p_1 \cdots p_r$ is the product of a power of 2 and distinct odd primes p_i where $p_i - 1$ is a power of 2. It is an elementary exercise to see that a prime p with $p - 1$ a power of 2 must be of the form

$$p = 2^{2^s} + 1$$

for some integer s . Such primes are called *Fermat primes*. The first few are

$$3 = 2^1 + 1$$

$$5 = 2^2 + 1$$

$$17 = 2^4 + 1$$

$$257 = 2^8 + 1$$

$$65537 = 2^{16} + 1$$

(but $2^{32} + 1$ is not a prime, being divisible by 641). It is not known if there are infinitely many Fermat primes. We summarize this in the following proposition.

Proposition 29. The regular n -gon can be constructed by straightedge and compass if and only if $n = 2^k p_1 \cdots p_r$ is the product of a power of 2 and distinct Fermat primes.

The proof above actually indicates a procedure for constructing the regular n -gon as a succession of square roots. For example, the construction of the regular 17-gon (solved by Gauss in 1796 at age 19) requires the construction of the subfields of degrees 2, 4, 8 and 16 in $\mathbb{Q}(\zeta_{17})$. These subfields can be constructed by forming the *periods* of ζ_{17} as in the example of the 13th roots of unity above. In this case, the fact that $\mathbb{Q}(\zeta_{17})$ is obtained by a series of quadratic extensions reflects itself in the fact that the periods can be “halved” successively (i.e., if $H_1 < H_2$ are subgroups with $[H_2 : H_1] = 2$ then the periods for H_1 satisfy a quadratic equation whose coefficients involve the periods for H_2). For example, the periods for the subgroup of index 2 (generated by σ_2) in the Galois group are ($\zeta = \zeta_{17}$)

$$\begin{aligned}\eta_1 &= \zeta + \zeta^2 + \zeta^4 + \zeta^8 + \zeta^9 + \zeta^{13} + \zeta^{15} + \zeta^{16} \\ \eta_2 &= \zeta^3 + \zeta^5 + \zeta^6 + \zeta^7 + \zeta^{10} + \zeta^{11} + \zeta^{12} + \zeta^{14}\end{aligned}$$

which “halve” the period for the full Galois group and which satisfy

$$\eta_1 + \eta_2 = -1$$

(from the minimal polynomial satisfied by ζ_{17}) and

$$\eta_1 \eta_2 = -4$$

(which requires computation — we know that it must be rational by Galois Theory, since this product is fixed by all the elements of the Galois group). Hence these two periods are the roots of the quadratic equation

$$x^2 + x - 4 = 0$$

which we can solve explicitly. In a similar way, the periods for the subgroup of index 4 (generated by σ_4) naturally halve these periods, so are quadratic over these, etc. In this way one can determine ζ_{17} explicitly in terms of iterated square roots. For example, one finds that $8(\zeta + \zeta^{-1}) = 16 \cos\left(\frac{2\pi}{17}\right)$ (which is enough to construct the regular 17-gon) is given explicitly by

$$-1 + \sqrt{17} + \sqrt{2(17 - \sqrt{17})} + 2\sqrt{17 + 3\sqrt{17} - \sqrt{2(17 - \sqrt{17})}} - 2\sqrt{2(17 + \sqrt{17})}.$$

A relatively simple construction of the regular 17-gon (shown to us by J.H. Conway) is indicated in the exercises.

While we have seen that it is not possible to solve for ζ_n using only successive square roots in general, by definition it is possible to obtain ζ_n by successive extraction of higher roots (namely, taking an n^{th} root of 1). This is not the case for solutions of general equations of degree n , where one cannot generally determine solutions by radicals, as we shall see in the next sections.

EXERCISES

1. Determine the minimal polynomials satisfied by the primitive generators given in the text for the subfields of $\mathbb{Q}(\zeta_{13})$.
2. Determine the subfields of $\mathbb{Q}(\zeta_8)$ generated by the periods of ζ_8 and in particular show that not every subfield has such a period as primitive element.
3. Determine the quadratic equation satisfied by the period $\alpha = \zeta_5 + \zeta_5^{-1}$ of the 5th root of unity ζ_5 . Determine the quadratic equation satisfied by ζ_5 over $\mathbb{Q}(\alpha)$ and use this to explicitly solve for the 5th root of unity.
4. Let $\sigma_a \in \text{Gal}(\mathbb{Q}(\zeta_n)/\mathbb{Q})$ denote the automorphism of the cyclotomic field of n^{th} roots of unity which maps ζ_n to ζ_n^a where a is relatively prime to n and ζ_n is a primitive n^{th} root of unity. Show that $\sigma_a(\zeta) = \zeta^a$ for every n^{th} root of unity.
5. Let p be a prime and let $\epsilon_1, \epsilon_2, \dots, \epsilon_{p-1}$ denote the primitive p^{th} roots of unity. Set $p_n = \epsilon_1^n + \epsilon_2^n + \dots + \epsilon_{p-1}^n$, the sum of the n^{th} powers of the ϵ_i . Prove that $p_n = -1$ if p does not divide n and that $p_n = p - 1$ if p does divide n . [One approach: $p_1 = -1$ from $\Phi_p(x)$; show that p_n is a Galois conjugate of p_1 for p not dividing n , hence is also -1 .]
6. Let ζ_n denote a primitive n^{th} root of unity and let $K = \mathbb{Q}(\zeta_n)$ be the associated cyclotomic field. Let a denote the trace of ζ_n from K to \mathbb{Q} (cf. Exercise 18 of Section 2). Prove that $a = 1$ if $n = 1$, $a = 0$ if n is divisible by the square of a prime, and $a = (-1)^r$ if n is the product of r distinct primes.
7. Show that complex conjugation restricts to the automorphism $\sigma_{-1} \in \text{Gal}(\mathbb{Q}(\zeta_n)/\mathbb{Q})$ of the cyclotomic field of n^{th} roots of unity. Show that the field $K^+ = \mathbb{Q}(\zeta_n + \zeta_n^{-1})$ is the subfield of real elements in $K = \mathbb{Q}(\zeta_n)$, called the *maximal real subfield of K*.
8. Let $K_n = \mathbb{Q}(\zeta_{2^{n+2}})$ be the cyclotomic field of 2^{n+2} -th roots of unity, $n \geq 0$. Set $\alpha_n = \zeta_{2^{n+2}} + \zeta_{2^{n+2}}^{-1}$ and $K_n^+ = \mathbb{Q}(\alpha_n)$, the maximal real subfield of K_n .
 - Show that for all $n \geq 0$, $[K_n : \mathbb{Q}] = 2^{n+1}$, $[K_n : K_n^+] = 2$, $[K_n^+ : \mathbb{Q}] = 2^n$, and $[K_{n+1}^+ : K_n^+] = 2$.
 - Determine the quadratic equation satisfied by $\zeta_{2^{n+2}}$ over K_n^+ in terms of α_n .
 - Show that for $n \geq 0$, $\alpha_{n+1}^2 = 2 + \alpha_n$ and hence show that

$$\alpha_n = \sqrt{2 + \sqrt{2 + \sqrt{\cdots + \sqrt{2}}}} \quad (\text{n times}),$$

giving an explicit formula for the (constructible) 2^{n+2} -th roots of unity.

9. Notation as in the previous exercise.
 - Prove that K_n^+ is a cyclic extension of \mathbb{Q} of degree 2^n . [Use an explicit isomorphism $(\mathbb{Z}/2^{n+2}\mathbb{Z})^\times \cong \mathbb{Z}/2\mathbb{Z} \times \mathbb{Z}/2^n\mathbb{Z}$ as abelian groups (i.e., $(\mathbb{Z}/2^{n+2}\mathbb{Z})^\times$ is isomorphic to a cyclic group of order 2 and a cyclic group of order 2^n — cf. Exercises 22 and 23 of Section 2.3)]
 - Prove that K_n is a biquadratic extension of K_{n-1}^+ and that two of the three intermediate subfields are K_n^+ and K_{n-1} . Prove that the remaining field intermediate between K_{n-1}^+ and K_n is a cyclic extension of \mathbb{Q} of degree 2^n .
10. Prove that $\mathbb{Q}(\sqrt[3]{2})$ is not a subfield of any cyclotomic field over \mathbb{Q} .
11. Prove that the primitive n^{th} roots of unity form a basis over \mathbb{Q} for the cyclotomic field of n^{th} roots of unity if and only if n is squarefree (i.e., n is not divisible by the square of any prime).

12. Let σ_p denote the Frobenius automorphism $x \mapsto x^p$ of the finite field \mathbb{F}_q of $q = p^n$ elements. Viewing \mathbb{F}_q as a vector space V of dimension n over \mathbb{F}_p we can consider σ_p as a linear transformation of V to V . Determine the characteristic polynomial of σ_p and prove that the linear transformation σ_p is diagonalizable over \mathbb{F}_p if and only if n divides $p - 1$, and is diagonalizable over the algebraic closure of \mathbb{F}_p if and only if $(n, p) = 1$.
13. Let $n = p_1^{a_1} p_2^{a_2} \cdots p_k^{a_k}$ be the prime factorization of n and let ζ_n be a primitive n^{th} root of unity. For each $i = 1, 2, \dots, k$ define d_i by $n = p_i^{a_i} d_i$ and let $\zeta_{p_i^{a_i}} = \zeta_n^{d_i}$, so that $\zeta_{p_i^{a_i}}$ is a particular primitive $p_i^{a_i}$ -th root of unity. Let $\sigma_a \in \text{Gal}(\mathbb{Q}(\zeta_n)/\mathbb{Q})$ be the automorphism mapping ζ_n to ζ_n^a for a relatively prime to n .
- (a) Prove that for $i = 1, 2, \dots, k$, σ_a maps $\zeta_{p_i^{a_i}}$ to $\zeta_{p_i^{a_i}}^a$ and gives an automorphism of $\mathbb{Q}(\zeta_{p_i^{a_i}})/\mathbb{Q}$ which depends only on $a \pmod{p_i^{a_i}}$, which we may denote $\sigma_{a \pmod{p_i^{a_i}}}$.
 - (b) Prove that the map $\sigma_a \mapsto (\sigma_{a \pmod{p_1^{a_1}}}, \dots, \sigma_{a \pmod{p_k^{a_k}}})$ is the isomorphism of Corollary 27 corresponding to the Chinese Remainder Theorem for $(\mathbb{Z}/n\mathbb{Z})^\times$.

The following Exercises 14 to 18 determine the periods associated to a primitive 17^{th} root of unity and provide a proof for the simple geometric construction indicated in Exercise 17 for the regular 17-gon. Let $\zeta = \zeta_{17} = \cos \frac{2\pi}{17} + i \sin \frac{2\pi}{17}$ be a fixed primitive 17^{th} root of unity in \mathbb{C} .

14. Define the *periods* of ζ as follows:

$$\begin{array}{ll} \eta_1 = \zeta + \zeta^2 + \zeta^4 + \zeta^8 + \zeta^9 + \zeta^{13} + \zeta^{15} + \zeta^{16} & \eta'_3 = \zeta^6 + \zeta^7 + \zeta^{10} + \zeta^{11} \\ \eta_2 = \zeta^3 + \zeta^5 + \zeta^6 + \zeta^7 + \zeta^{10} + \zeta^{11} + \zeta^{12} + \zeta^{14} & \eta'_4 = \zeta^3 + \zeta^5 + \zeta^{12} + \zeta^{14} \\ \eta'_1 = \zeta + \zeta^4 + \zeta^{13} + \zeta^{16} & \eta''_1 = \zeta + \zeta^{16} \\ \eta'_2 = \zeta^2 + \zeta^8 + \zeta^9 + \zeta^{15} & \eta''_2 = \zeta^4 + \zeta^{13}. \end{array}$$

- (a) Show that all of these periods are real numbers and that $\eta''_1 = 2 \cos \frac{2\pi}{17}$. Show that as real numbers these periods are approximately

$$\begin{array}{lll} \eta_1 \sim 1.562 & \eta'_1 \sim 2.049 & \eta'_3 \sim -2.906 \\ \eta_2 \sim -2.562 & \eta'_2 \sim -0.488 & \eta'_4 \sim 0.344 \\ & & \eta''_1 \sim 1.865 \\ & & \eta''_2 \sim 0.185. \end{array}$$

- (b) Prove that η_1 and η_2 are roots of the equation $x^2 + x - 4 = 0$.
(c) Prove that η'_1 and η'_2 are roots of the equation $x^2 - \eta_1 x - 1 = 0$ and that η'_3 and η'_4 are roots of the equation $x^2 - \eta_2 x - 1 = 0$.
(d) Prove that η''_1 and η''_2 are roots of the equation $x^2 - \eta'_1 x + \eta'_4 = 0$.

15. Prove that if $\tan 2\theta = a$ ($0 < 2\theta < \frac{\pi}{2}$) then $\tan \theta$ satisfies the equation $x^2 - \frac{2}{a}x - 1 = 0$.
16. Let C be the circle in \mathbb{R}^2 having the points (h, k) and $(0, 1)$ as a diameter. Prove that this circle intersects the x -axis if and only if $h^2 - 4k \geq 0$ and in this case the two intercepts are the roots of the equation $x^2 - hx + k = 0$.
17. (*Construction of the Regular 17-gon*) Draw a circle of radius 2 centered at the origin $(0, 0)$.
- (a) Join the point $(4, 0)$ to the point $(0, 1)$ and construct the line ℓ_1 bisecting the angle

between this line and the y -axis. Construct the line ℓ_2 perpendicular to ℓ_1 in Figure 2.

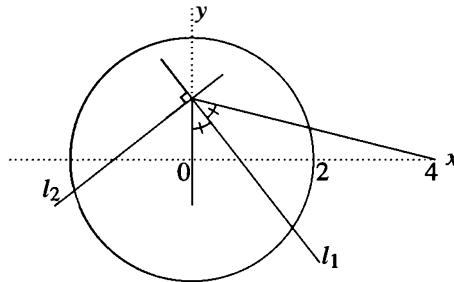


Fig. 2

- (b) Using the intersection of ℓ_1 and the x -axis as center and radius equal to the distance to $(0, 1)$, construct the circle C_1 and let $A = (s, 0)$ be the right-hand point of intersection of C_1 with the x -axis. Similarly, let $B = (t, 0)$ denote the right-hand point of intersection of the x -axis and the circle C_2 whose center is the intersection of ℓ_2 and the x -axis and whose radius is equal to the distance to $(0, 1)$ as in Figure 3.

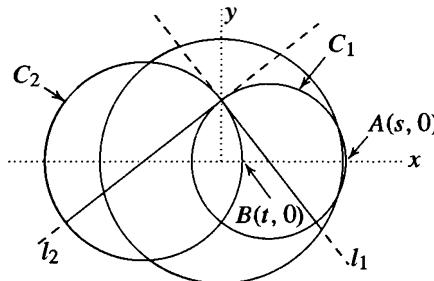


Fig. 3

- (c) Construct a perpendicular to the x -axis at the point A and mark off the distance t from $(0, 0)$ to B to construct the point (s, t) . Construct the circle with (s, t) and $(0, 1)$ as a diameter and let P denote the right-hand point of intersection of this circle with the x -axis. The perpendicular to the x -axis at P intersects the circle of radius 2 at the second vertex of a regular 17-gon whose first vertex is at $(2, 0)$, hence constructs the regular 17-gon by straightedge and compass as in Figure 4.

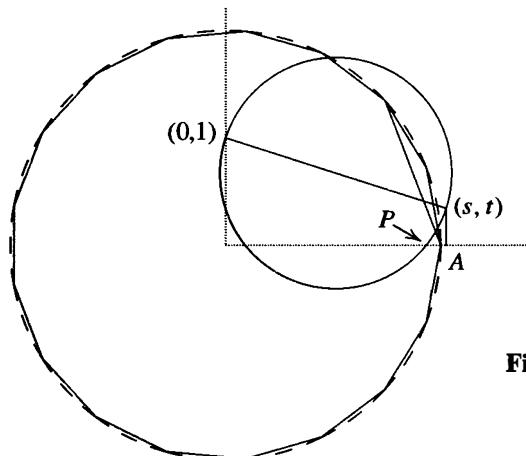


Fig. 4