

1. Algebraic grammars

~~we have studied~~
Let M be a monoid and let Ξ be
(disjoint from M).
a finite alphabet. We define a monoid
 $M[\Xi]$ as follows. The elements of $M[\Xi]$
are words.

(1.1) $m_0 \xi_1 m_1 \dots m_{k-1} \xi_k m_k$
with $k \geq 0$, $m_0, \dots, m_k \in M$ and $\xi_1, \dots, \xi_k \in \Xi$.
The elements m_0, \dots, m_k are called the
coefficients while ξ_1, \dots, ξ_k are called the
variables. The product of (1.1) with a word
~~marked~~

$n_0 \xi_1 n_1 \dots n_{l-1} \xi_l n_l$
is the ~~associated~~ word

$m_0 \xi_1 m_1 \dots m_{k-1} \xi_k (m_k n_0) \xi_1 n_1 \dots n_{l-1} \xi_l n_l$
We shall say that the monoid $M[\Xi]$
is obtained from M by adjoining the
set Ξ as a set of free variables.

If $M = \Sigma^*$ then $M[\Xi] = (\Sigma \cup \Xi)^*$.

~~assuming that Σ and Ξ are disjoint.~~

An algebraic (or context-free)

grammar, G , in a monoid M consists of the following data:

- (1.2) A finite (non-empty) alphabet Σ disjoint from M . The elements of Σ are called variables.
- (1.3) A distinguished element $s_0 \in \Sigma$ called the start variable.
- (1.4) A finite subset R of $\Sigma \times M[\Sigma]$ called the set of rules.

For each rule $r \in R$, we denote by $\bar{r} \in \Sigma$ its first coordinate and by $r \in M[\Sigma]$ its second coordinate.

Instead of writing $(\bar{r}, r) \in R$ we shall frequently write $\bar{r} \rightarrow r$ or $r: \bar{r} \rightarrow r$

and $x \notin \Sigma$. Thus

$$x \in (\Sigma \cup \Xi)^+ - \underline{\underline{\Sigma}}$$

In other words, an positive grammar has
no rules of the form

$$\underline{\underline{S}} \rightarrow T, \quad \underline{\underline{S}} \rightarrow \underline{\underline{S'}}.$$

For each rule $r \in R$ we define a partial function $\cdot r : M[\Xi] \rightarrow M[\Xi]$
as follows:

$$w r = w'$$

whenever

$$w = u \bar{r} v, \quad u \in M[\Xi], \quad v \in M$$

$$w' = u r v$$

In all other cases $w r = \emptyset$. We call the reader's attention to the condition $v \in M$. It signifies that \bar{r} is rightmost appearance of a variable in w . It is only then that \bar{r} is replaced by r (i.e. is rewritten using the rule r).

Given

$$d = r_1 \dots r_n \in R^*$$

we define

$$wd = (\dots ((wr_1)r_2) \dots r_n)$$

and $w1 = w$ (if $n=0$). Thus each $d \in R$ defines a partial function $M[\Sigma] \rightarrow M[\Sigma]$ and

$$w(d_1 d_2) = (wd_1)d_2$$

If $wd = w'$ we say that d is a derivation of w into w' and also use the notation $d : w \rightarrow w'$ or $w \xrightarrow{d} w'$.

So far the start variable s_0 played no role. It does play a role in the following definition.

The language L ~~of G~~ defined by the grammar G , is the subset of Σ^*

$$(1.5) \quad L = \{ s \mid s \in \Sigma^*, d : s_0 \rightarrow s \text{ for some } d \in R^* \}$$

Since $s_0 \notin M$ we may replace $d \in R^*$ by

So far the start variable played no role. It does play a role in the following definition. A derivation

$$d: \xi_0 \rightarrow m, m \in M$$

is called a successful derivation. The set L of all elements $m \in M$ for which such a successful derivation exists is called the language determined by the grammar G and is denoted by GL . Thus

$$L = \{m \mid m \in M, d: \xi_0 \rightarrow m \text{ for some } d \in R^*\}$$

Since $\xi_0 \notin M$, we may replace $d \in R^*$ by $d \in R^+$.

The above definition does not take into account the number of successful derivations $d: \xi_0 \rightarrow m$ that may exist for a particular m . Thus the definition of GL as given above ignores multiplicity.

cies and defines GL as an ordinary or as B-subset of M.

If we wish to pay attention to the phenomenon of multiplicity, then we define mL , to be the number of successful derivations $d : \xi_0 \rightarrow m$. Since this number could be infinite we have

$$mL \in \mathbb{N} = N \cup \{\infty\}$$

where N is the set of all integers $n \geq 0$.

Thus L becomes a function

$$(L : M \rightarrow \mathbb{N}$$

or in other words an \mathbb{N} -subset of M.

An \mathbb{N} -subset or a B-subset of M is called algebraic if it is defined by some algebraic grammar in M.

We denote the class of all algebraic \mathbb{N} -subsets of M by

MAlg or

while the class of all algebraic B-subsets
is denoted by

MAlg
 δN .

Proposition 1.1. Every derivation

$$d : w_2 w_1 \rightarrow m, m \in M$$

admits a unique factorization

$$d = d_1 d_2, \quad d_1 : w_1 \rightarrow m_1, \quad d_2 : w_2 \rightarrow m_2$$

such that

$$(1.5) \quad d_1 : w_1 \rightarrow m_1, \quad d_2 : w_2 \rightarrow m_2, \quad m_2 m_1 = m$$

Conversely given (1.5) we have

$$d_1 d_2 : w_2 w_1 \rightarrow m_2 m_1.$$

Proof. We first show the uniqueness of the factorization $d = d_1 d_2$. Indeed assume that $d = e_1 e_2$ is another such factorization.

Assume $|d_1| < |e_1|$. Then $e_1 = d_1 c$. Since $w_1 d_1 = m \in M$ we have $w_1 e_1 = mc$. This is impossible since $mc = \emptyset$.

The proof of the existence of the factorization $d = d_1 d_2$ proceeded by induction on $|d|$. If $|d| = 0$, then

$w_1 w_2 = m$ and $d_1 = 1 = d_2$. Now assume that $|d| = l + n$ and that the assertion holds for derivations of length n .

We also observe that if $w_1 \in M$ then

$(w_2 w_1)d = (w_2 d)w_1$, and thus $d_1 = d$ and
 $d_2 = 1$ give the required factorization.

Let $d = rd'$. Then $w_2 w_1 = u\bar{r}p$ with
 $p \in M$. Since we assume that $w_1 \notin M$ it

follows that $w_1 = v\bar{r}p$, $u = w_2 v$. Then

$w_2 w_1 r = w_2 v \bar{r} p$. Since $d' : w_2(v\bar{r}p) \rightarrow m$

we have a factorization $d' = d'_1 d'_2$ such

that $(v\bar{r}p)d_1 = m_1$, $w_2 d_2 = m_2$ and $m = m_2 m_1$.

Then $d = (r d'_1) d'_2$ is the required factoriza-
tion of d ■

Exercise 1.1. Show that for each algebraic
K-subset ^A_M of M (with K = B or K = JY) there
exists a finitely generated submonoid M' of
M such that A is an algebraic K-subset
of M' .

2. Examples

Example 2.1. Let A be a finite subset of monoid M . Consider the grammar with a single variable ξ (which is then necessarily the start variable) and with rules

$$\xi \rightarrow a$$

one for each $a \in A$. Clearly A is the language ~~can~~ defined by this grammar.

Thus A is an algebraic subset of M .

This example also shows that if we allow the set R of rules to be infinite then any subset of M could be shown to be algebraic (see however Theorem V).

Example 2.2. Let $M = \{\sigma, \tau\}^*$. Consider the grammar

$$\xi \rightarrow \sigma \xi \tau, \quad \xi \rightarrow 1$$

Since there is only one variable ξ , it is automatically the start variable. Let r_1 and r_2 designate the two displayed rules. The

only derivations $d: \xi \rightarrow w$ possible are

$$r_1^n: \xi \rightarrow \sigma^n \xi \tau^n, \quad n \geq 0$$

$$r_2 r_1^n: \xi \rightarrow \sigma^n \tau^n, \quad n \geq 0$$

The derivations $r_2 r_1^n$ are the only successful ones and consequently the grammar defines the algebraic set

$$A = \{\sigma^n \tau^n \mid n \geq 0\}$$

If in this example we replace M by the commutative monoid freely generated by σ and τ , we obtain the subset $A = (\sigma \tau)^*$.

13

Example 2.3. Let $M = \{\sigma, \tau\}^*$, Consider

the grammar with variables ξ (start variable) and η as

$$r_1: \xi \rightarrow \sigma \xi \tau, r_2: \xi \rightarrow \tau \eta \sigma, r_3: \xi \rightarrow 1$$

$$r_4: \eta \rightarrow \tau \eta \sigma, r_5: \eta \rightarrow 1$$

Show that, the grammar computes the subset

$$A = \{ \sigma^p \tau^q \sigma^q \tau^p \mid p \geq 0, q \geq 0 \}$$

Show that the only successful derivations are

$$r_1^{p_1} r_2^{q_1} r_4^{q-1} r_5 : \xi \rightarrow \sigma^p \tau^q \sigma^q \tau^p \text{ if } q > 0$$

$$r_1^p r_3^q : \xi \rightarrow \sigma^p \tau^p \text{ if } q = 0$$

Example 2.4. Given any alphabet Σ consider grammar with rules

$$r_0: \xi \rightarrow \sigma \xi \sigma, r_0': \xi \rightarrow \sigma, r_0'': \xi \rightarrow 1$$

Show that this grammar defines the set of all palindromes

$$A = \{ s \mid s \in \Sigma^*, s = s^s \}$$

where s^s is the reversal of s .

Example 2.5. Let $M = \sigma^*$. Consider the grammar

$$\xi_0 \rightarrow \sigma, \xi \rightarrow \xi', \xi' \rightarrow \sigma$$

The grammar defines the set consisting of σ alone. However note that we have two distinct derivations $d: \xi_0 \rightarrow \sigma$, one of length 1 and another one of length 2. Thus the N-subset A of σ^* defined by this grammar is $\{ \sigma \}$.

15

3. Weights.

Let K be a commutative semiring, and let M be a monoid. An (M, K) -grammar G is an algebraic grammar (Σ, S_0, R) in M together with a weight function

$$\mu : R \rightarrow K$$

Thus each rule r has a weight $r\mu \in K$. This function μ is extended to R^* by setting

$$d\mu = (r_1\mu) \cdots (r_n\mu)$$

for $d = r_1 \cdots r_n$ (Note that $d\mu = 1$ if $n=0$).

Using these weights we now modify the definition of the language L defined by the grammar. L is now to become a K -subset of M i.e. a function

$$L : M \rightarrow K$$

defined by

$$(3.1) \quad m L = \sum_{\xi_0 d = m} d\mu$$

3.1

There is something drastically wrong
with this definition. The summation in (3.1)
extends over all derivations.

$$(3.2) \quad d : \xi_0 \rightarrow m$$

and there may be infinitely many of them.

Thus the sum (3.1) may be undefined
if K is just a semiring.

One way to remedy this difficulty
is to use a complete semiring K
rather than just a semiring (see A,
In such a semiring arbitrary summations
can be carried out and thus (3.1) is
well defined. The semirings \mathbb{B} and \mathbb{W}
are examples of complete semirings.

Actually the assumption that K is
complete is too strong. In III.3 we

2010

shall introduce the notion of a continuous semiring. In such a semiring countable sums can be performed and thus suffices to give meaning to (3.1). The class of algebraic K -subsets of M thus defined will be denoted by

$$(3.3) \quad M\text{Alg}_K$$

There is a second important case in which the summation (3.1) is legitimate. If we replace M by a Semigroup S which is loose (see next section for a definition) and if we assume that the grammar is positive (see next section for a definition), then it can be shown that for each $s \in S$, the derivations $d : \xi_0 \rightarrow s$ form a finite set. Thus (3.1) is again well defined, K . The class of

(for any commutative semiring)

positive

algebraic K -subsets of S thus defined
will be denoted by

$$(3.4) \quad S \text{Alg}_K$$

There is no conflict between the notation
(3.3) and (3.4) since in (3.3) M is a
monoid, while in (3.4) S is a loose
semigroup and such a semigroup cannot
have a unit element.

There is a slight conflict of notation
involving $M\text{Alg}_B$ and $M\text{Alg}_N$ since these
we introduced in section 1 without weights
(or, what is equivalent, with all ~~weights~~
the rules given weight 1). In the case
of $M\text{Alg}_B$ the conflict is resolved
trivially. Since $B = \{0, 1\}$ all weights are
0 or 1. Rules of weight 0 may be
discarded so that only rules of weight 1
remain. In the case of $M\text{Alg}_N$ the

conflict is resolved by

Proposition 3.1. For any (M, \mathcal{N}) -gramm. G there exists an (M, \mathcal{N}) grammar G' with all rules having weight 1 and such that G and G' define the same \mathcal{N} -subset of M .

Proof. Let $G = (\Sigma, S_0, R)$ and let $\mu: R \rightarrow \mathcal{N}$ be the weight function. The grammar is modified in the following manner: For each rule $r: \xi \rightarrow w$ in R the following four cases are considered

- (i) $r\mu = 0$. The rule r is discarded
- (ii) $r\mu = 1$. The rule r is left unchanged
- (iii) $r\mu = k$ with $1 < k < \infty$. New variables η_1, \dots, η_k are introduced and the rule r is replaced by the $2k$ rules

$$\xi \rightarrow \eta_i, \quad \eta_i \rightarrow w, \quad i=1, \dots, k.$$
- (iv) $r\mu = \infty$. A new variable η is introduced and the rule r is replaced by

P
the three rules

$$\xi \rightarrow \gamma, \gamma \rightarrow \eta, \eta \rightarrow w$$

All the new rules introduced in (iii) and (iv) receive the weight 1. The verification that the new grammar defines the same Σ -subset of M as G is immediate \blacksquare

The following notational device is useful in handling grammars with weights. If $d: \xi \rightarrow w$ is a derivation with weight $d\mu = k$, then we shall write $d: \xi \rightarrow kw$. In particular $r: \xi \rightarrow kw$ signifies that the rule r has weight k .

4. Loose semigroups and positive grammars

A semigroup S is said to be loose if each $s \in S$ admits only a finite number of factorizations

$$(4.1) \quad s = s_1 \dots s_n, \quad s_1, \dots, s_n \in S$$

It follows that a loose semigroup S cannot contain an idempotent. Thus S cannot have a unit element or a zero element and cannot contain a non-empty finite subsemigroup.

If we denote by $|s|$ the largest integer n for which a factorization (4.1) exists, we obtain a function $S \rightarrow \mathbb{N}$ with the following properties

$$(4.2) \quad |s| > 0 \text{ for all } s \in S$$

$$(4.3) \quad |s| = 1 \text{ iff } s \text{ is indecomposable,}$$

i.e. iff $s \in S - S^2$

$$(4.4) \quad |st| \geq |s| + |t|.$$

2.2

Let Σ be a set equipped with a bijection $\varphi: \Sigma \rightarrow S - S^2$. This bijection is extended to a morphism $\varphi: \Sigma^+ \rightarrow S$. Since S is generated by its indecomposable elements it follows that the morphism φ is surjective.

Let \sim be the congruence relation in Σ^+ defined by $w \sim w'$ iff $w\varphi = w'\varphi$.

The following two properties of this congruence are easily verified

(4.5) Each $\sigma \in \Sigma$ is a congruence class in itself

(4.6) Each congruence class is finite.

Conversely give any congruence \sim in Σ^+ satisfying (4.5) and (4.6), the quotient semigroup $S = \Sigma^+ / \sim$ is loose. This shows the extent by which loose semigroups differ from free ones.

25

Example 4.1. Each free semigroup Σ^+ is
loose. Each free commutative semigroup Σ_c^+

is loose.

(Σ_c^+ is a groupoid with identity and no
inverses.)

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

Each free semigroup Σ^+ is loose.

Each free commutative semigroup Σ_c^+ is loose.

L. 6

Given any loose semigroup S , we denote by S^* the monoid obtained from S by adjoining a unit element. Let K be any commutative semiring. An (S, K) -grammar is defined to be an (S^*, K) -grammar $G = (\Sigma, \delta_S, R)$ which has no rules $\xi \rightarrow 1$. If further G has no rules $\xi \rightarrow \xi'$ then we say that G is a positive (S, K) -grammar.

We define δ

$$S[\Sigma] = S^*[\Sigma] - 1$$

We note that $S[\Sigma]$ is a loose semigroup and that for any word

$$w = m_0 \xi_{i_0} m_1 \dots m_{p-1} \xi_{i_p} m_p$$

we have

$$|w| = p + |m_0| + \dots + |m_p|$$

where $|m_i| = 0$ whenever $m_i = 1 \in S^*$. The integer p is called the degree of w and is denoted

$$ww^* \delta = w\delta + w^*\delta$$

$$\|ww\| \geq \|w\| + \|w'\| ?$$

$$\|w\| = \|x \xi y\|$$

$$\geq \|x\| + \|\xi\| + \|y\| \geq 1 + \|x\| + \|y\|$$

$$= 1 + \|x \xi y\|$$

by ws.

The fact that the grammar is positive may be expressed by the inequality

$$(4.7) \quad 2 \leq 2|x| - x\delta$$

Indeed, if $x=1$ then $|x|=0=x\delta$
and if $x=\xi$ then $|x|=x\delta=1$. In
all other cases the inequality (4.7) holds.

We claim that in a positive grammar
the inequality

$$(4.8) \quad \|d\| + 2\|u\| - u\delta \leq 2\|w\| - w\delta$$

holds for any derivation

$$(4.9) \quad d: u \rightarrow \not w$$

The inequality clearly holds if $d=1$ and $u=w$. Assume that $d=d'\circ r$ and that the inequality holds for d' . Thus d is the composition

$$u \xrightarrow{d} x\xi y \xrightarrow{r} x\underline{r} y = w$$

$$\therefore \|d\| + \|w\| \leq \|wd\|$$

$$\|d\| + \|u\| = 1 + \|d'\| + \|u\delta\| \not\leq \|w\|$$

$$\leq 1 + \|ud'\| \leq \|ud'\|$$

Then

$$\begin{aligned}
 |d| + 2|w| - w\delta &= | + d' + 2|w| - w\delta \\
 &\leq | + 2|x\circ y| - (x\circ y)\delta \\
 &= 2 + 2|x| + 2|y| - x\delta - y\delta \\
 &\leq 2 + 2|\bar{x}y| - 2|\bar{x}| + (\bar{x}\bar{y})\delta - \bar{x}\delta \\
 &= 2|w| - w\delta - (2|\bar{x}| - \bar{x}\delta - 2) \\
 &\leq 2|w| - w\delta
 \end{aligned}$$

as required.

Taking $w = \xi$ and $w = s \in S$ in (4.9)

we obtain

Proposition 4.1. If

$$d : \xi \rightarrow s, s \in S$$

\Rightarrow a derivation in positive (S, K) -grammar,

then

$$|d| < 2|s| \blacksquare$$

Corollary 4.2. For each $s \in S$, there is
only a finite number of derivations. \therefore

$$d : \xi \rightarrow s \blacksquare$$

Exercise 4.1. Let $\Sigma = \{\sigma\}$ be a single letter alphabet and let $A = 2\sigma$ be the N -subset of Σ^+ consisting of the letter σ with multiplicity 2. Show that A is positive algebraic, but that every positive grammar $G = (\Sigma, S_0, R)$ defining A must contain a rule $S_0 \rightarrow \sigma$ with weight 2 and multiplicity 2.

Compare this with Proposition 3.1.

Exercise 4.2. For each of the sets A in

Examples 2.2 - 2.3 construct a positive grammar for the set A with unit element removed.

Exercise 4.3. Show that if $\varphi: S \rightarrow S'$ is a morphism of loose semigroups then $|s| \leq |s\varphi|$ for all $s \in S$.

5. Linguistic interpretation; parsing

Temporary graph
pp. 28-34

6. Introduction to the fixed point approach.

In Chapters II and III we shall develop a completely different approach to algebraic sets, completely sidestepping derivations. To prepare the reader, we shall give here a heuristic discussion of a special case ~~that will nevertheless~~ which however is sufficient to give the reader a taste of the things to come.

H. B. Johnstone

We return to the grammar

$$\xi \rightarrow \sigma \xi \tau , \quad \xi \rightarrow 1$$

considered in Example I, 1.2 and defining
the subset

$$A = \{ \sigma^n \tau^n \mid n \geq 0 \}$$

of $M = \{\sigma, \tau\}^*$.

We shall interpret ξ as representing
an unknown subset X of M . The
rule $\xi \rightarrow \sigma \xi \tau$ will, ^{thus} be interpreted as
the inclusion $\sigma X \tau \subset X$ and similarly
the rule $\xi \rightarrow 1$ will be interpreted
as $1 \in X$. Using + instead of \cup
we thus obtain the inclusion

$$(1.1) \quad \dagger X \tau + 1 \subset X$$

The set A does satisfy this inclusion
but much more is true. Indeed let F
be any subset of M satisfying (1.1)

Then $1 \in B$ and since $\sigma B \tau \subset B$, it follows inductively that $\sigma^n \tau^n \in B$ for all $n \geq 0$. Thus $A \subset B$, and B is

the minimal solution of the inclusion. Next we note that for A we actually have

$$A' = \sigma A \tau + 1$$

so that A also is the minimal solution of the equation.

$$X = \sigma X \tau + 1$$

The process of passing from grammars to equations is not limited to the example above. The grammar

$$\begin{aligned} S &\rightarrow \sigma \gamma, \quad S \rightarrow S^2 \tau, \quad S \rightarrow \sigma \\ \gamma &\rightarrow S \gamma, \quad \gamma \rightarrow \tau^2, \quad \gamma \rightarrow 1 \end{aligned}$$

in which x, y are variables and
 $M = \{r, t\}^*$ leads to the system
of equations

$$X = rY + X^2t + r$$

$$Y = XY + t^2 + 1$$

On the right hand side of each equation we have a polynomial
(in the non-commutative sense) with
 X and Y as variables and with
elements of M as coefficients. Jointly
the two equations may be written as

$$(6.2) \quad (X, Y) = (X, Y)P$$

where $P : 2^M \times 2^M \rightarrow 2^M \times 2^M$

$$\therefore P : 2^M \times 2^M \rightarrow 2^M \times 2^M$$

is a polynomial transformation in
which the various monomial correspond
to the rules of the grammar.

positive

In Chapter II we shall treat (S, K) -grammars where S is a loose semigroup and K is any commutative semiring. In this case it will be shown that the equation (6.2) has a unique solution, and that the first coordinate of this solution is the language defined by the grammar. This is the positive case.

In Chapter III we shall define continuous semirings and study (M, K) -grammars where M is any monoid and K is a continuous semiring. The equation (6.2) may then have many solutions, but one of them is the smallest (in the sense of a partial order inherent in K). The language defined by the grammar is then the first coordinate of the minimal solu-

tion of (6.2). This is the continuous
case.