

A Coalgebraic View on Context-Free Languages and Streams

Joost Winter

Centrum Wiskunde & Informatica

May 10, 2011

Overview

- ▶ The coalgebraic picture of regular languages and expressions, and likewise that of rational streams and power series, is well-known.
- ▶ It is interesting to see how this work can be extended to, in first instance, context-free languages.
- ▶ In this presentation, a coalgebraic treatment of context-free languages through systems of behavioural differential equations is given.
- ▶ This definition format can be generalized to arbitrary formal power series (in noncommuting variables), including streams, yielding a notion of *context-free power series* and *streams*.
- ▶ Some examples of streams that are found to be 'context-free' in this sense are given.

Formal power series, formal languages, and streams

- ▶ Given a finite set A called the *alphabet*, and a semiring R , a *formal power series* on A with coefficients in R is a function

$$A^* \rightarrow R.$$

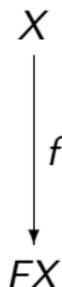
- ▶ When R is the Boolean semiring $\{0, 1\}$ (with $1 + 1 = 1$), formal power series on A with coefficients in R correspond to formal languages over the alphabet A .
- ▶ When A is a singleton set, formal power series on A with coefficients on R correspond to *streams* over R .

F-Coalgebras

Given a functor F , an F -coalgebra consists of a tuple (X, f) :

- ▶ X is a set, the *carrier set*.
- ▶ f is a function from X to FX .

Diagrammatically:



Coalgebras representing formal power series (1)

In the this talk, we will be concerned with coalgebras over functors of the type $R \times (-)^A$. But what does this mean?

- ▶ R is some semiring.
- ▶ \times is the cartesian product.
- ▶ A is a finite set called the *alphabet*.
- ▶ $(-)$ is a placeholder for the carrier set.
- ▶ X^A denotes the function space from A to X .

Coalgebras representing formal power series (2)

So: a coalgebra (X, f) over the functor $R \times (-)^A$ consists of a set X and a function f that maps every $x \in X$ to an element $f(x) \in R \times X^A$.

In this talk, we will use the following notation:

- ▶ Given $x \in X$, $o(x)$ (called the *output value of x*) will be the first component of $f(x)$.
- ▶ Given $x \in X$, x_a (called the *a -derivative of x*) will be the second component of $f(x)$, applied to a .

So: for every x , $o(x)$ is an element of the semiring R , and for every x and a , x_a is an element of X again.

When (in the case of streams) the alphabet is a singleton set $\{a\}$, we usually write x' instead of x_a .

Coalgebras representing formal power series (3)

We can extend the notion of derivatives from alphabet symbols (i.e. elements of A) to words (i.e. elements of A^*) inductively:

- ▶ $x_\lambda = x$
- ▶ $x_{a \cdot w} = (x_a)_w$.

Homomorphisms and bisimulations (1)

Given two $R \times (-)^A$ -coalgebras (X, f) and (Y, g) , a function $h : X \rightarrow Y$ is a *homomorphism* if the following hold:

1. For every $x \in X$, $o(x) = o(h(x))$.
2. For every $x \in X$ and $a \in A$, $h(x_a) = (h(x))_a$.

Homomorphisms and bisimulations (2)

Given two $R \times (-)^A$ -coalgebras (X, f) and (Y, g) , a relation $R \subseteq X \times Y$ is a *bisimulation* if the following hold:

1. If $(x, y) \in R$, then $o(x) = o(y)$.
2. If $(x, y) \in R$, then for all $a \in A$, $(x_a, y_a) \in R$.

Final coalgebras (1)

Consider the $2 \times (-)^A$ -coalgebra (\mathcal{L}, l) defined as follows:

- ▶ \mathcal{L} is the set of all languages on the alphabet A .
- ▶ For any $L \in \mathcal{L}$:
 - ▶ $o(L)$ is 1 iff the empty word is in L .
 - ▶ $F_a = \{w \mid a \cdot w \in L\}$.

This is a *final coalgebra*: for every $2 \times (-)^A$ -coalgebra (X, f) , there is a *unique* homomorphism h from (X, f) to (\mathcal{L}, l) .

Given a $2 \times (-)^A$ -coalgebra (X, f) , and an element $x \in X$, we let $\llbracket x \rrbracket$ denote the value of x under this unique homomorphism.

Final coalgebras (2)

Generalizing the picture from the previous slide to arbitrary semirings R , the final coalgebras will be sets of formal power series:

- ▶ \mathcal{S} is the set of all formal power series on A with coefficients in R , i.e. the function space from A^* to R .
- ▶ For any $F \in \mathcal{S}$:
 - ▶ $o(F) = F(\lambda)$.
 - ▶ $F_a = G : G(x) = F(a \cdot x)$.

Again, we let $\llbracket x \rrbracket$ denote the value of x under the final homomorphism.

Languages and streams

The presentation given above is a generalization of both coalgebras representing languages and coalgebras representing streams:

- ▶ When we choose the Boolean semiring $\{0, 1\}$, we get coalgebras representing languages over the alphabet A .
- ▶ When we choose a singleton alphabet A , we get coalgebras representing streams over the semiring R .

The next few slides will be about coalgebras for the functor $\mathcal{D} := 2 \times (-)^A$ of formal languages.

The coalgebra of regular expressions

The set \mathcal{E} of *regular expressions* over a finite alphabet A and the semiring $(2, +, \cdot, 0, 1)$ can be defined as follows:

$$t ::= a \in A \mid r \in 2 \mid t + t \mid t \cdot t \mid t^*$$

We can assign a \mathcal{D} -coalgebra structure to this set of regular expressions by specifying the output values and derivatives for each expression, giving us a \mathcal{D} -coalgebra (\mathcal{E}, e) :

t	$o(t)$	t_a
$r \in 2$	r	0
$b \in A$	0	if $b = a$ then 1 else 0
$u + v$	$o(u) + o(v)$	$u_a + v_a$
$u \cdot v$	$o(u) \cdot o(v)$	$u_a \cdot v + o(u) \cdot v_a$
u^*	1	$u_a \cdot u^*$

Kleene's theorem, coalgebraically

- ▶ For any language $L \in \mathcal{L}$, we define the subcoalgebra generated by \mathcal{L} as

$$\langle L \rangle := \{L_w \mid w \in A^*\}$$

It is easy to see that this indeed generates a subcoalgebra: given any $K \in \langle L \rangle$, it is easy to see that for every $a \in A$, also $K_a \in \langle L \rangle$. In other words, $\langle L \rangle$ is closed under taking derivatives to alphabet symbols.

- ▶ Kleene's theorem, coalgebraically (Rutten, 1998): *For any $L \in \mathcal{L}$, $\langle L \rangle$ is finite iff there is a regular expression t such that $L = \llbracket t \rrbracket$.*

Introduction: context-free grammars and languages

- ▶ The 'next step up' from regular expressions and languages, and finite automata, in the Chomsky hierarchy, are the context-free languages and grammars, and pushdown automata.
- ▶ We will present a format of coinductively defined systems of equations: it turns out that these systems of equations characterize precisely the context-free languages.

Systems of equations (1)

We will use terms t specified as follows:

$$t ::= a \in A \mid x \in X \mid r \in 2 \mid t + t \mid t \cdot t$$

where X is a finite set of variables, and A , as before, is a finite alphabet. Given X , we let TX denote the set of terms over X . A well-formed system of equations, for a set of variables X , consists of:

1. For every $x \in X$, exactly one equation of the form $o(x) = v$, where $v \in \{0, 1\}$.
2. For every $x \in X$ and $a \in A$, exactly one equation of the form $x_a = t$, where $t \in TX$.

Systems of equations (2)

Alternatively, we can consider a well-formed system of equations as a mapping

$$f : X \rightarrow 2 \times TX^A$$

We can extend such a mapping f to the \mathcal{D} -coalgebra (TX, \bar{f}) generated by (X, f) as follows:

t	$o(t)$	t_a
$x \in X$	$o(x)$	x_a (as specified by f)
$r \in 2$	r	0
$b \in A$	0	if $b = a$ then 1 else 0
$u + v$	$o(u) + o(v)$	$u_a + v_a$
$u \cdot v$	$o(u) \cdot o(v)$	$u_a \cdot v + o(u) \cdot v_a$

Systems of equations (3)

- ▶ This construction can be summarized diagrammatically:

$$\begin{array}{ccccc}
 X & \hookrightarrow & TX & \xrightarrow{\quad} & \mathcal{L} \\
 \downarrow f & & \swarrow \bar{f} & & \downarrow l \\
 & & & & 2 \times \mathcal{L}^A \\
 & & & \xrightarrow{\quad} & \\
 2 \times TX^A & & & &
 \end{array}$$

$\begin{array}{c} \parallel \\ \parallel \end{array}$

- ▶ Proposition: A language L is context-free iff there is a well-formed system of equations (X, f) and an $x \in X$, such that $\llbracket x \rrbracket = L$ w.r.t. the coalgebra (TX, \bar{f}) generated by it.

From CFGs to systems of equations (1)

- ▶ We say a context-free grammar is in weak Greibach normal form, if every production rule has a right hand side either equal to the empty word λ , or of the form $a \cdot t$.
- ▶ As the name implies, this is a weakening of the more familiar Greibach normal form. As a direct result, every CFG can be represented in weak Greibach normal form.

From CFGs to systems of equations (2)

We transform a CFG G in weak Greibach normal form into a system of equations as follows:

- ▶ We let the set X of variables be equal to the set of nonterminals in the grammar.
- ▶ Given a $x \in X$, we set $o(x) = 1$ iff the grammar contains a production rule $x \rightarrow \lambda$.
- ▶ Given a $x \in X$ and an $a \in A$, we set

$$x_a = \sum \{w \mid x \rightarrow a \cdot w\}$$

Given an initial symbol $x_0 \in X$, we now have $o((x_0)_w) = 1$ (and, hence, $w \in \llbracket x_0 \rrbracket$) iff w is in the language generated by G .

From systems of equations to CFGs

Conversely, given a system of equations, we can construct a CFG in weak Greibach normal form:

- ▶ We first transform the system of equations to a new, equivalent, system, in which all derivatives are in disjunctive normal form, and do not contain any superfluous 0s or 1s.
- ▶ Derivatives in this new system are disjunctions of sequences of alphabet symbols and variables.
- ▶ We let the grammar include a rule $x \rightarrow \lambda$ whenever $o(x) = 1$.
- ▶ We let the grammar include a rule $x \rightarrow a \cdot w$, whenever w is a sequence of alphabet symbols and variables occurring as a disjunct in x_a .

Generalizing context-freeness

- ▶ Note that most of the definitions on the earlier slides do not necessarily require the underlying semiring to be the Boolean semiring!
- ▶ Taking the definition of *regular expressions* and applying it to arbitrary semirings, we obtain rational streams and power series, which are well-known.
- ▶ Furthermore, we can easily generalize the notion above, of *context-free languages*, to *context-free streams* and *context-free power series*.

Context-free streams (1)

When we take the semiring of natural numbers as underlying semiring, the Catalan numbers

1, 1, 2, 5, 14, 42, 132, 429, 1430, . . .

are context-free, and are generated by the following system of equations:

$$o(x) = 1 \quad x' = x \cdot x$$

Context-free streams (2)

When we take the Boolean *field* (with $1 + 1 = 0$) as underlying semiring, the Prouhet-Thue-Morse sequence

1001011001101001...

is context-free, and is generated by the following system of equations:

$$\begin{aligned}o(x) &= 1 & x' &= y \\o(y) &= 0 & y' &= z \\o(z) &= 0 & z' &= x + y + z + z \cdot w \\o(w) &= 1 & w' &= y \cdot w\end{aligned}$$

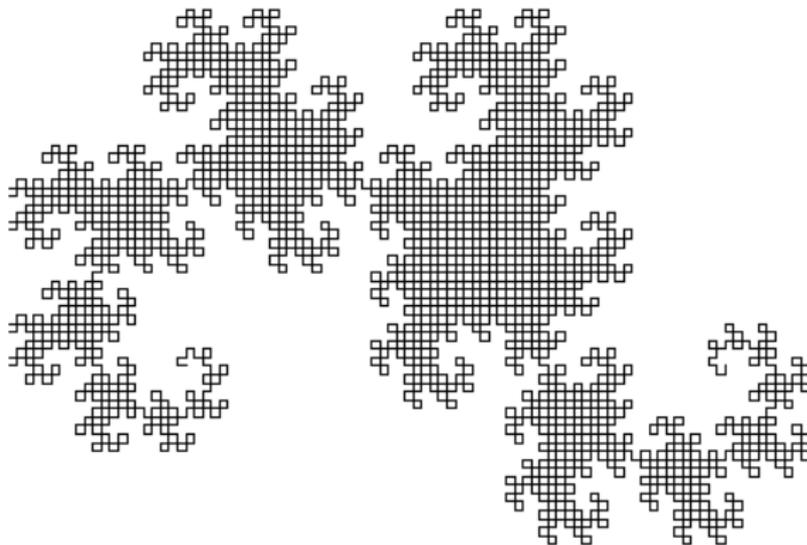
Context-free streams (3)

Still using the Boolean field as underlying semiring, the context-free system of equations

$$\begin{array}{ll} o(x) = 1 & x' = y \\ o(y) = 1 & y' = z \\ o(z) = 0 & z' = w \\ o(w) = 1 & w' = v \\ o(v) = 1 & v' = v + w + v \cdot x + x \cdot x \end{array}$$

gives us the paperfolding sequence, or Dragon curve sequence. . .

Context-free streams (4)



Conclusions and further work

- ▶ There is a very neat coalgebraic representation of regular expressions, and Kleene's theorem can be expressed succinctly in a coalgebraic fashion.
- ▶ We have extended this work towards context-free languages and grammars, and provided a coalgebraic characterization using systems of equations.
- ▶ The first steps towards a generalization to other functors of the type $R \times (-)^A$ has been made, and there are some neat examples of context-free streams.
- ▶ Future work: further investigate this notion of 'generalized context-freeness' of power series, and see how this relates to other, existing notions.

Bibliography

-  [Jacobs/Rutten, 1997] Bart Jacobs, Jan Rutten, *A Tutorial on (Co)Algebras and (Co)Induction*
-  [Rutten, 1998] Jan Rutten, *Automata and Coinduction (An Exercise in Coalgebra)*
-  [Rutten, 2005] Jan Rutten, *A Coinductive Calculus of Streams*
-  [Silva, 2010] Alexandra Silva, *Kleene Coalgebra*
-  [Winter/Bonsangue/Rutten, 2011] Joost Winter, Marcello Bonsangue, Jan Rutten, *Context-free Languages, Coalgebraically*