

Language-Processing Problems

Roland Backhouse
DIMACS, 8th July, 2003

Introduction

“Factors” and the “factor matrix” were introduced by Conway (1971).

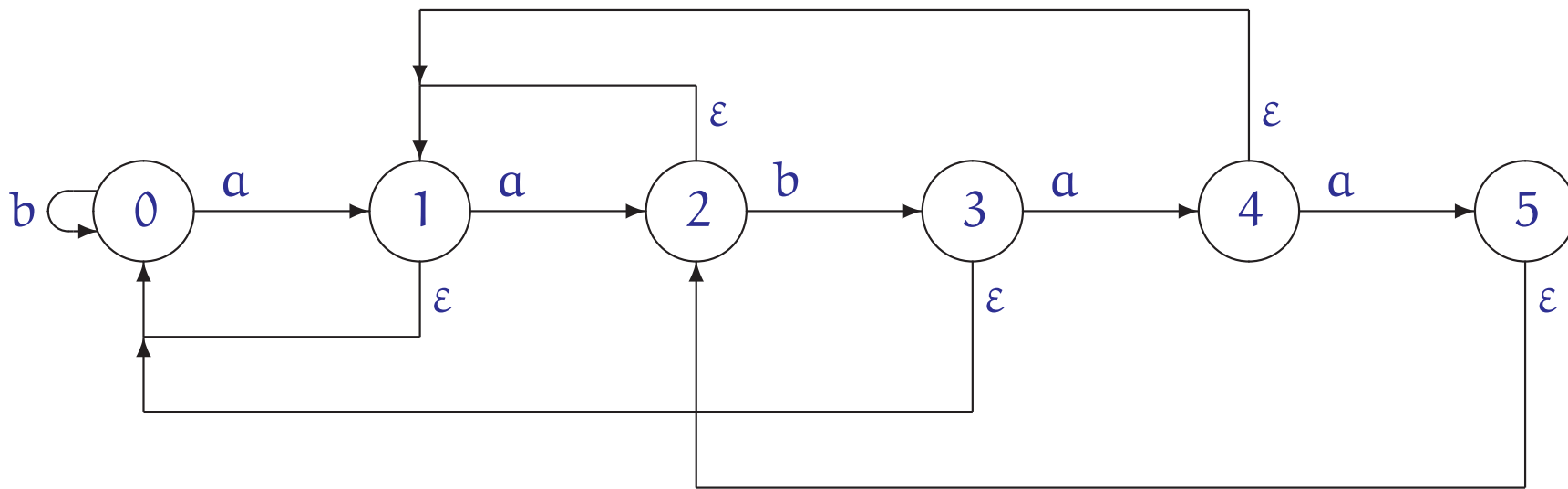
He used them very effectively in, for example, constructing biregulators.

Conway’s discussion is wordy, making it difficult to understand. There are also occasional errors which are difficult to detect and add to the confusion. (“The theorem does prevent **E** from occurring twice” should read “The theorem does *not* prevent **E** from occurring twice.”)

KMP Failure Function (pattern aabaa)

node i	1	2	3	4	5
failure node $f(i)$	0	1	0	1	2

Factor Graph (language $\Sigma^* aabaa$)



Language Problems

$$S ::= aSS \mid \varepsilon .$$

Is-empty

$$S = \phi \equiv (\{a\} = \phi \vee S = \phi \vee S = \phi) \wedge \{\varepsilon\} = \phi .$$

Nullable

$$\varepsilon \in S \equiv (\varepsilon \in \{a\} \wedge \varepsilon \in S \wedge \varepsilon \in S) \vee \varepsilon \in \{\varepsilon\} .$$

Shortest word length

$$\#S = (\#a + \#S + \#S) \downarrow \#\varepsilon .$$

Non-Example

$$aa \in S \not\equiv (aa \in \{a\} \wedge aa \in S \wedge aa \in S) \vee aa \in \{\varepsilon\} .$$

Fusion

Many problems are expressed in the form

$$\textit{evaluate} \circ \textit{generate}$$

where *generate* generates a (possibly infinite) candidate set of solutions, and *evaluate* selects a best solution.

Examples:

$$\textit{shortest} \circ \textit{path} ,$$

$$(\text{x}\in) \circ \text{L} .$$

Solution method is to *fuse* the generation and evaluation processes, eliminating the need to generate all candidate solutions.

Conditions for Fusion

Fusion is made possible when

- *evaluate* is an adjoint in a *Galois connection*,
- *generate* is expressed as a *fixed point*.

Algorithms for solving resulting fixed point equation include

- brute-force iteration,
- Knuth's generalisation of Dijkstra's shortest path algorithm. .

Solution method typically involves *generalising* the problem.

Galois Connections

Suppose $\mathcal{A} = (A, \sqsubseteq)$ and $\mathcal{B} = (B, \preceq)$ are partially ordered sets and suppose $F \in A \leftarrow B$ and $G \in B \leftarrow A$. Then (F, G) is a *Galois connection of \mathcal{A} and \mathcal{B}* iff, for all $x \in B$ and $y \in A$,

$$F(x) \sqsubseteq y \equiv x \preceq G(y) .$$

Examples

Negation:

$$\neg p \Rightarrow q \equiv p \Leftarrow \neg q .$$

Ceiling function:

$$\lceil x \rceil \leq n \equiv x \leq n .$$

Maximum:

$$x \uparrow y \leq z \equiv x \leq z \wedge y \leq z .$$

Even (divisible by two):

$$\text{if } b \rightarrow 2 \square \neg b \rightarrow 1 \text{ fi } \setminus m \equiv b \Rightarrow \text{even}(m) .$$

Parsing

$$x \in S \Rightarrow b \quad \equiv \quad S \subseteq \text{if } b \rightarrow \Sigma^* \square \neg b \rightarrow \Sigma^* - \{x\} \text{ fi} .$$

Shortest Word (Path)

Let $\Sigma^{\geq k}$ denote the set of all words over alphabet Σ of length at least k .

Let $\#S$ denote the length of a shortest word in the language S .

$$\#S \geq k \quad \equiv \quad S \subseteq \Sigma^{\geq k} .$$

(Most common application is when S is the set of paths from one node to another in a graph.)

Fusion Theorem

$$F(\mu_{\preceq} g) = \mu_{\sqsubseteq} h$$

provided that

- F is a lower adjoint in a Galois connection of \sqsubseteq and \preceq (see brief summary of definition below)
- $F \circ g = h \circ F$.

Galois Connection

$$F(x) \sqsubseteq y \equiv x \preceq G(y) \text{ .}$$

F is called the *lower* adjoint and G the *upper* adjoint.

Language Recognition

Problem: For given word x and grammar G , determine $x \in L(G)$.

That is, implement

$$(x \in) \circ L \ .$$

Language $L(G)$ is the least fixed point (with respect to the subset relation) of a monotonic function.

$(x \in)$ is the lower adjoint in a Galois connection of languages (ordered by the subset relation) and booleans (ordered by implication).

(Recall,

$$x \in S \Rightarrow b \quad \equiv \quad S \subseteq \text{if } b \rightarrow \Sigma^* \square \neg b \rightarrow \Sigma^* - \{x\} \text{ fi} \ .)$$

Nullable Languages

Problem: For given grammar G , determine $\varepsilon \in L(G)$.

$$(\varepsilon \in) \circ L$$

Solution: Easily expressed as a fixed point computation.

Works because:

- The function $(\mathbf{x} \in)$ is a lower adjoint in a Galois connection (for all \mathbf{x} , but in particular for $\mathbf{x} = \varepsilon$).
- For all languages S and T ,

$$\varepsilon \in S \cdot T \quad \equiv \quad \varepsilon \in S \wedge \varepsilon \in T \quad .$$

Problem Generalisation

Problem: For given grammar G , determine whether all words in $L(G)$ have even length. I.e. implement

$$\text{alleven} \circ L \text{ .}$$

The function **alleven** is a lower adjoint in a Galois connection. Specifically, for all languages S and T ,

$$\text{alleven}(S) \Leftarrow b \quad \equiv \quad S \subseteq \text{if } \neg b \rightarrow \Sigma^* \square b \rightarrow (\Sigma \cdot \Sigma)^* \text{ fi} \text{ .}$$

Nevertheless, fusion *doesn't* work (directly) because

- there is no \otimes such that, for all languages S and T ,

$$\text{alleven}(S \cdot T) \quad \equiv \quad \text{alleven}(S) \otimes \text{alleven}(T) \text{ .}$$

Solution: Generalise by tupling: compute simultaneously **alleven** and **allodd**.

General Context-Free Parsing

Problem: For given grammar G , determine $x \in L(G)$.

$$(x \in) \circ L \text{ .}$$

Not (in general) expressible as a fixed point computation.

Fusion *fails* because: for all x , $x \neq \varepsilon$, there is no \otimes such that, for all languages S and T ,

$$x \in S \cdot T \equiv (x \in S) \otimes (x \in T) \text{ .}$$

CYK: Let $F(S)$ denote the relation $\langle i, j :: x[i..j] \in S \rangle$.

Works because:

- The function F is a lower adjoint.
- For all languages S and T ,

$$F(S \cdot T) = F(S) \bullet F(T)$$

where $B \bullet C$ denotes the composition of relations B and C .

Language Inclusion

Problem: For fixed (regular) language E and varying S , determine

$$S \subseteq E .$$

Example: Emptiness test:

$$S \subseteq \phi .$$

Example: Pattern Matching: given pattern P , for each prefix t of text T , evaluate:

$$\{t\} \subseteq \Sigma^* \cdot \{P\} .$$

Example: All words are of even length:

$$S \subseteq (\Sigma \cdot \Sigma)^* .$$

Language Inclusion

Problem: For fixed (regular) language E and varying S , determine

$$S \subseteq E .$$

- Function $(\subseteq E)$ is a lower adjoint. Specifically,

$$S \subseteq E \Leftarrow b \quad \equiv \quad S \subseteq \text{if } b \rightarrow E \square \neg b \rightarrow \Sigma^* \text{ fi} .$$

- But, for $E \neq \phi$ and $E \neq \Sigma^*$, there is no \otimes such that, for all languages S and T ,

$$S \cdot T \subseteq E \quad \equiv \quad (S \subseteq E) \otimes (T \subseteq E) .$$

Solution (Oege de Moor): Use factor theory to derive generalisation.

Factors

For all languages S , T and U ,

$$S \cdot T \subseteq U \equiv T \subseteq S \setminus U ,$$

$$S \cdot T \subseteq U \equiv S \subseteq U / T .$$

Note:

$$S \setminus (U / T) = (S \setminus U) / T .$$

Hence, write

$$S \setminus U / T .$$

Left and Right Factors

Define the functions \triangleleft and \triangleright by

$$\begin{aligned} X_{\triangleleft} &= E/X , \\ X_{\triangleright} &= X \setminus E . \end{aligned}$$

By definition, the range of \triangleleft is the set of *left* factors of E and the range of \triangleright is the set of *right* factors of E .

We also have the Galois connection:

$$X \subseteq Y_{\triangleleft} \equiv Y \subseteq X_{\triangleright} .$$

Hence,

$$\begin{aligned} X_{\triangleleft\triangleright\triangleleft} &= X_{\triangleleft} , \\ X_{\triangleright\triangleleft\triangleright} &= X_{\triangleright} , \\ E_{\triangleleft\triangleright} &= E = E_{\triangleright\triangleleft} . \end{aligned}$$

The Factor Matrix

Let \mathcal{L} denote the set of left factors of E .

Define the *factor matrix* of E to be the binary operator \setminus restricted to $\mathcal{L} \times \mathcal{L}$. Thus entries in the matrix take the form $L_0 \setminus L_1$ where L_0 and L_1 are left factors of E .

The factor matrix of E is denoted by $\llbracket E \rrbracket$. It is a reflexive, transitive matrix.

$$\llbracket E \rrbracket = \llbracket E \rrbracket^* .$$

The row and column containing individual factors, the left factors, the right factors, and E itself, is given by:

$$U \setminus E / V = U_{\triangleright\triangleleft} \setminus V_{\triangleleft} ,$$

$$V_{\triangleleft} = E_{\triangleleft} \setminus V_{\triangleleft} ,$$

$$U_{\triangleright} = U_{\triangleright\triangleleft} \setminus E_{\triangleright\triangleleft} ,$$

$$E = E_{\triangleleft} \setminus E_{\triangleright\triangleleft} .$$

Using the Factor Matrix

Problem: For fixed regular language E and varying S , determine

$$S \subseteq E .$$

Generalisation: For fixed regular language E and varying S , determine the relation

$$S \subseteq \llbracket E \rrbracket .$$

(Formally, the relation $\langle L, M :: S \subseteq L \setminus M \rangle$ where L and M range over the left factors of E .)

Works because:

$$S \cdot T \subseteq \llbracket E \rrbracket \quad \equiv \quad (S \subseteq \llbracket E \rrbracket) \bullet (T \subseteq \llbracket E \rrbracket) .$$

where $B \bullet C$ denotes the composition of relations B and C .

Proof

We have to show that

$$S \cdot T \subseteq U \triangleleft W \equiv \langle \exists V :: S \subseteq U \triangleleft V \wedge T \subseteq V \triangleleft W \rangle .$$

First,

$$\begin{aligned} & S \cdot T \subseteq E \\ = & \{ \text{unit of conjunction} \} \\ & S \cdot T \subseteq E \wedge \text{true} \\ = & \{ \text{factors, } T \triangleleft = E/T; \text{ cancellation} \} \\ & S \subseteq T \triangleleft \wedge T \triangleleft \cdot T \subseteq E \\ = & \{ \text{factors, } T \triangleleft \triangleright = T \triangleleft \setminus E \} \\ & S \subseteq T \triangleleft \wedge T \subseteq T \triangleleft \triangleright . \end{aligned}$$

Whence:

$$\begin{aligned}
& S \cdot T \subseteq U_{\triangleleft} \setminus W_{\triangleleft} \\
= & \quad \{ \text{factors, definition of } W_{\triangleleft} \} \\
& U_{\triangleleft} \cdot S \cdot T \cdot W \subseteq E \\
= & \quad \{ \text{above, with } S, T := U_{\triangleleft} \cdot S, T \cdot W \} \\
& U_{\triangleleft} \cdot S \subseteq (T \cdot W)_{\triangleleft} \wedge T \cdot W \subseteq (T \cdot W)_{\triangleleft} \\
= & \quad \{ \text{factors} \} \\
& S \subseteq U_{\triangleleft} \setminus (T \cdot W)_{\triangleleft} \wedge T \subseteq (T \cdot W)_{\triangleleft} / W \\
= & \quad \{ U_{\triangleright} / W = U \setminus W_{\triangleleft} \} \\
& S \subseteq U_{\triangleleft} \setminus (T \cdot W)_{\triangleleft} \wedge T \subseteq (T \cdot W)_{\triangleleft} \setminus W_{\triangleleft} . \\
\Rightarrow & \quad \{ \text{one-point rule} \} \\
& \langle \exists V :: S \subseteq U_{\triangleleft} \setminus V_{\triangleleft} \wedge T \subseteq V_{\triangleleft} \setminus W_{\triangleleft} \rangle \\
\Rightarrow & \quad \{ \text{Leibniz} \} \\
& \langle \exists V :: S \cdot T \subseteq U_{\triangleleft} \setminus V_{\triangleleft} \cdot V_{\triangleleft} \setminus W_{\triangleleft} \rangle \\
\Rightarrow & \quad \{ \text{cancellation,} \} \\
& S \cdot T \subseteq U_{\triangleleft} \setminus W_{\triangleleft} .
\end{aligned}$$

Summary

- Use of **fusion** as programming method.
- **Problem generalisation** involves generalising the **algebra** in the solution domain.
- **Factor theory** as basis for language inclusion problems.

Challenges

- Efficient computation of factor matrices.
- Extension to non-regular languages.

References

J.H. Conway, “Regular Algebra and Finite Machines”, Chapman and Hall, London, 1971.

Backhouse, R.C and Carré, B.A. “Regular algebra applied to path-finding problems”, J. Institute of Mathematics and its Applications, vol. 15, pp. 161–186, 1975.

Backhouse, R.C. and Lutz, R.K., “Factor graphs, failure functions and bi-trees”, Automata, Languages and Programming, LNCS 52, pp. 61–75, 1977.

Roland Backhouse, “Fusion on Languages”, 10th European Symposium on Programming, LNCS 2028, pp. 107–121, 2001.

O. de Moor, S. Drape, D. Lacey and G. Sittampalam. Incremental program analysis via language factors.

(www.comlab.ox.ac.uk/oucl/work/oege.demoor/pubs.htm)

For related publications on fixed points, Galois connections and mathematics of program construction, see

www.cs.nott.ac.uk/~rcb/papers