

FUNCTIONAL PEARL

A program to solve Sudoku

RICHARD BIRD

Programming Research Group, Oxford University
Wolfson Building, Parks Road, Oxford OX1 3QD, UK
email: bird@comlab.ox.ac.uk

There's no maths involved. You solve
the puzzle with reasoning and logic.
Advice on how to play Sudoku,
The Independent Newspaper

1 Introduction

The game of Sudoku is played on a 9 by 9 board. Given is a matrix of characters, such as

2	1	.	3	8
.	5
.	7	.	.	.	6	.	.	.
.	1	3
.	9	8	1	.	.	2	5	7
3	1	8	.	.
9	.	.	8	.	.	.	2	.
.	5	.	.	6	9	7	8	4
4	.	.	2	5

The idea is to fill in the dots with the digits 1 to 9 so that each row, each column and each of the component 3 by 3 boxes contains the digits 1 to 9 exactly once. In general there may be one, none or many solutions, though in a good Sudoku puzzle there is always a unique solution. Our aim in this pearl is to derive a Haskell program to solve Sudoku puzzles. Specifically, we will define a function

$$\text{solve} :: \text{Board} \rightarrow [\text{Board}]$$

for computing all the ways a given board may be filled. If we want only one solution we can take the head of the list. Lazy evaluation means that only the first result will then be computed.

We do not want our program to depend on the board size, as long as it is of the form $(N^2 \times N^2)$ for positive N , nor on the precise characters chosen for the entries. Instead, the program is parameterized by three constants, *boardsize*, *boxsize* and *cellvals*, and one test *blank* $:: \text{Char} \rightarrow \text{Bool}$ for determining whether a given entry

is blank. For concreteness, we can take

$$\begin{aligned} \text{boardsize} &= 9 \\ \text{boxsize} &= 3 \\ \text{cellvals} &= \text{“123456789”} \\ \text{blank} &= (= \text{‘ ’}) \end{aligned}$$

Changing cell values, e.g. to “TONYBLAIR” is easy.

2 Specification

The first aim is to write down the simplest and clearest specification of Sudoku without regard to how efficient the result might be. Such a specification will help us focus ideas on how a more efficient solution might be obtained, as well as being a starting point for program manipulation.

One possibility is first to construct a list of all correctly completed boards, and then to test the given board against these boards to identify those whose entries match the given ones. Another possibility, and the one we will take, is to start with the given board and to generate all possible completions. Each completed board is then tested to see if it is correct, that is, does not contain duplicate entries in each row, column or box.

A board is a matrix of characters:

$$\begin{aligned} \text{type Matrix } a &= [[a]] \\ \text{type Board} &= \text{Matrix Char} \end{aligned}$$

Strictly speaking a given board should first be checked to see that every non-blank entry is an element of *cellvals*. Invalid boards should be rejected. However, for simplicity we will assume that the given board does satisfy the basic requirements.

The function *correct* tests whether a filled board, that is, one containing no blank characters, has different entries in each row, column and box:

$$\begin{aligned} \text{correct} &:: \text{Board} \rightarrow \text{Bool} \\ \text{correct } b &= \text{all nodups (rows } b) \wedge \\ &\quad \text{all nodups (cols } b) \wedge \\ &\quad \text{all nodups (boxs } b) \end{aligned}$$

The function *nodups* can be defined by

$$\begin{aligned} \text{nodups} &:: \text{Eq } a \Rightarrow [a] \rightarrow \text{Bool} \\ \text{nodups } [] &= \text{True} \\ \text{nodups } (x : xs) &= \text{notElem } x \text{ } xs \wedge \text{nodups } xs \end{aligned}$$

2.1 Rows, columns and boxes

If a matrix is given by a list of its rows, the function *rows* is just the identity function on matrices:

$$\begin{aligned} \text{rows} &:: \text{Matrix } a \rightarrow \text{Matrix } a \\ \text{rows} &= \text{id} \end{aligned}$$

We have, trivially, that $\text{rows} \cdot \text{rows} = \text{id}$.

The function *cols* computes the transpose of a matrix. One possible definition is:

$$\begin{aligned} \text{cols} &:: \text{Matrix } a \rightarrow \text{Matrix } a \\ \text{cols } [xs] &= [[x] \mid x \leftarrow xs] \\ \text{cols } (xs : xss) &= \text{zipWith } (:) xs (\text{cols } xss) \end{aligned}$$

We also have $\text{cols} \cdot \text{cols} = \text{id}$.

The boxes of a matrix can be computed by:

$$\begin{aligned} \text{boxs} &:: \text{Matrix } a \rightarrow \text{Matrix } a \\ \text{boxs} &= \text{map } \text{ungroup} \cdot \text{ungroup} \cdot \text{map } \text{cols} \cdot \text{group} \cdot \text{map } \text{group} \end{aligned}$$

The function *group* groups a list into component lists of length *boxsize*, and *ungroup* takes a grouped list and ungroups it:

$$\begin{aligned} \text{group} &:: [a] \rightarrow [[a]] \\ \text{group} &= \text{groupBy } \text{boxsize} \\ \text{ungroup} &:: [[a]] \rightarrow [a] \\ \text{ungroup} &= \text{concat} \end{aligned}$$

We omit the definition of *groupBy*. Using $\text{ungroup} \cdot \text{group} = \text{id}$ and $\text{group} \cdot \text{ungroup} = \text{id}$, it is easy to show that $\text{boxs} \cdot \text{boxs} = \text{id}$ by simple equational reasoning.

2.2 Generating choices and matrix cartesian product

The function *choices* replaces blank entries in a board with all possible choices for that entry. Using *Choices* as a synonym for *[Char]*, we have

$$\begin{aligned} \text{choices} &:: \text{Board} \rightarrow \text{Matrix } \text{Choices} \\ \text{choices} &= \text{map } (\text{map } \text{choose}) \\ \text{choose } e &= \text{if blank } e \text{ then cellvals else } [e] \end{aligned}$$

Of course, not every possible choice is valid for each cell, and we will return to this point later on.

The function *mcp* (matrix cartesian product) generates a list of all possible boards from a given matrix of choices:

$$\begin{aligned} \text{mcp} &:: \text{Matrix } [a] \rightarrow [\text{Matrix } a] \\ \text{mcp} &= \text{cp} \cdot \text{map } \text{cp} \end{aligned}$$

The function *cp* computes the cartesian product of a list of lists:

$$\begin{aligned} \text{cp} &:: [[a]] \rightarrow [[a]] \\ \text{cp } [] &= [[]] \\ \text{cp } (xs : xss) &= [x : ys \mid x \leftarrow xs, ys \leftarrow \text{cp } xss] \end{aligned}$$

Note that $\text{cp } xss$ returns an empty list if *xss* contains an empty list. Thus $\text{mcp } cm$ returns an empty list if any entry of *cm* is the empty list.

2.3 Specification

The function *sudoku* can now be defined by

$$\begin{aligned} \textit{sudoku} &:: \textit{Board} \rightarrow [\textit{Board}] \\ \textit{sudoku} &= \textit{filter correct} \cdot \textit{mcp} \cdot \textit{choices} \end{aligned}$$

However, this specification is executable in principle only. Assuming about a half of the 81 entries are fixed initially, there are about 9^{40} , or

$$147808829414345923316083210206383297601$$

boards to check! We therefore need a better approach.

3 Pruning the choices

Obviously, not every possible choice is valid for each cell. A better choice for a blank entry in row r , column c and box b is any cell value that does not appear among the fixed entries in row r , column c or box b . An entry in a matrix of choices is fixed if it is a singleton list. The fixed entries in a given row, column or box, are given by

$$\begin{aligned} \textit{fixed} &:: [\textit{Choices}] \rightarrow \textit{Choices} \\ \textit{fixed} &= \textit{concat} \cdot \textit{filter single} \end{aligned}$$

where $\textit{single} :: [a] \rightarrow \textit{Bool}$ tests whether the argument is a singleton list. The fixed entries can be removed from a list of choices by

$$\begin{aligned} \textit{reduce} &:: [\textit{Choices}] \rightarrow [\textit{Choices}] \\ \textit{reduce css} &= \textit{map} (\textit{remove} (\textit{fixed css})) \textit{css} \\ \textit{remove fs cs} &= \textbf{if } \textit{single cs} \textbf{ then } \textit{cs} \textbf{ else } \textit{delete fs cs} \end{aligned}$$

We leave the definition of *delete* to the reader.

Now, how shall we prune the matrix of choices? The aim is to define a function

$$\textit{prune} :: \textit{Matrix Choices} \rightarrow \textit{Matrix Choices}$$

satisfying the equation

$$\textit{filter correct} \cdot \textit{mcp} = \textit{filter correct} \cdot \textit{mcp} \cdot \textit{prune}$$

The function *prune* removes the fixed choices from each row, column or box. The question is a good test of one's programming ability for it seems easy to get into a mess. So, it is worthwhile adding a small pause at this point, to see if the reader can come up with a short definition that meets the requirement.

(Pause)

The calculational programmer would calculate a definition, and that is precisely what we are going to do.

The first step is to rewrite *filter correct* in the form

$$\begin{aligned} \textit{filter correct} &= \textit{filter} (\textit{all nodups} \cdot \textit{boxs}) \cdot \\ &\quad \textit{filter} (\textit{all nodups} \cdot \textit{cols}) \cdot \\ &\quad \textit{filter} (\textit{all nodups} \cdot \textit{rows}) \end{aligned}$$

The order of the component filters is unimportant.

Now we send these filters one by one into battle with *mcp*. We will need some weapons, the first of which is the law

$$\underline{\text{filter } (p \cdot f) = \text{map } f \cdot \text{filter } p \cdot \text{map } f}$$

which is valid provided $f \cdot f = \text{id}$. In particular, the law is valid if f is one of *rows*, *cols*, or *boxs*.

Another useful law is the following one:

$$\text{filter } (\text{all } p) \cdot \text{cp} = \text{cp} \cdot \text{map } (\text{filter } p) \quad \text{先枚举再过滤} = \text{先过滤再枚举}$$

In words, if we want only those lists all of whose elements satisfy p from a cartesian product, then we can obtain them by taking the cartesian product of the elements satisfying p of the component lists.

We will also need the following facts:

$$\begin{aligned} \text{map rows} \cdot \text{mcp} &= \text{mcp} \cdot \text{rows} \\ \text{map cols} \cdot \text{mcp} &= \text{mcp} \cdot \text{cols} \\ \text{map boxs} \cdot \text{mcp} &= \text{mcp} \cdot \text{boxs} \end{aligned}$$

These laws are intuitively clear and we will not verify them formally.

We need one final law: the crucial property of the function *reduce* defined above is that

$$\text{filter nodups} \cdot \text{cp} = \text{filter nodups} \cdot \text{cp} \cdot \text{reduce}$$

Here is the calculation. Let f be one of *rows*, *cols* or *boxs*:

$$\begin{aligned} & \text{filter } (\text{all nodups} \cdot f) \cdot \text{mcp} \\ = & \quad \{\text{since } \text{filter } (p \cdot f) = \text{map } f \cdot \text{filter } p \cdot \text{map } f \text{ if } f \cdot f = \text{id}\} \\ & \text{map } f \cdot \text{filter } (\text{all nodups}) \cdot \text{map } f \cdot \text{mcp} \\ = & \quad \{\text{since } \text{map } f \cdot \text{mcp} = \text{mcp} \cdot f \text{ if } f \in \{\text{boxs}, \text{cols}, \text{rows}\}\} \\ & \text{map } f \cdot \text{filter } (\text{all nodups}) \cdot \text{mcp} \cdot f \\ = & \quad \{\text{definition of mcp}\} \\ & \text{map } f \cdot \text{filter } (\text{all nodups}) \cdot \text{cp} \cdot \text{map } \text{cp} \cdot f \\ = & \quad \{\text{since } \text{filter } (\text{all } p) \cdot \text{cp} = \text{cp} \cdot \text{map } (\text{filter } p)\} \\ & \text{map } f \cdot \text{cp} \cdot \text{map } (\text{filter nodups} \cdot \text{cp}) \cdot f \\ = & \quad \{\text{property of reduce}\} \\ & \text{map } f \cdot \text{cp} \cdot \text{map } (\text{filter nodups} \cdot \text{cp} \cdot \text{reduce}) \cdot f \\ = & \quad \{\text{since } \text{filter } (\text{all } p) \cdot \text{cp} = \text{cp} \cdot \text{map } (\text{filter } p)\} \\ & \text{map } f \cdot \text{filter } (\text{all nodups}) \cdot \text{cp} \cdot \text{map } (\text{cp} \cdot \text{reduce}) \cdot f \\ = & \quad \{\text{since } \text{map } f \cdot \text{filter } p = \text{filter } (p \cdot f) \cdot \text{map } f \text{ if } f \cdot f = \text{id}\} \\ & \text{filter } (\text{all nodups} \cdot f) \cdot \text{map } f \cdot \text{mcp} \cdot \text{map } \text{reduce} \cdot f \\ = & \quad \{\text{since } \text{map } f \cdot \text{mcp} = \text{mcp} \cdot f \text{ if } f \in \{\text{boxs}, \text{cols}, \text{rows}\}\} \\ & \text{filter } (\text{all nodups} \cdot f) \cdot \text{mcp} \cdot f \cdot \text{map } \text{reduce} \cdot f \\ = & \quad \{\text{definition of pruneBy } f; \text{ see below}\} \\ & \text{filter } (\text{all nodups} \cdot f) \cdot \text{mcp} \cdot \text{pruneBy } f \end{aligned}$$

The definition of *pruneBy* is

$$\begin{aligned} \text{pruneBy} &:: (\text{MatrixChoices} \rightarrow \text{MatrixChoices}) \rightarrow \\ &\quad (\text{MatrixChoices} \rightarrow \text{MatrixChoices}) \\ \text{pruneBy } f &= f \cdot \text{map reduce} \cdot f \end{aligned}$$

We have shown that, provided *f* is one of *rows*, *cols* or *boxs*,

$$\text{filter } (\text{all nodups} \cdot f) \cdot \text{mcp} = \text{filter } (\text{all nodups} \cdot f) \cdot \text{mcp} \cdot \text{pruneBy } f$$

For the final step we need one more law, the fact that we can interchange the order of two *filter* operations:

$$\text{filter } p \cdot \text{filter } q = \text{filter } q \cdot \text{filter } p$$

This law is not generally valid in Haskell without qualification on the boolean functions *p* and *q*, but provided *p* and *q* are total functions, as is the case here, the law is OK. Indeed we implicitly made use of it when claiming that the order of the component filters in the expansion of *filter correct* was unimportant.

Now we can calculate, abbreviating *nodups* to *nd* to keep the expressions short:

$$\begin{aligned} &\text{filter correct} \cdot \text{mcp} \\ = &\quad \{\text{rewriting } \text{filter correct} \text{ as three filters}\} \\ &\text{filter } (\text{all nd} \cdot \text{boxs}) \cdot \text{filter } (\text{all nd} \cdot \text{cols}) \cdot \text{filter } (\text{all nd} \cdot \text{rows}) \cdot \text{mcp} \\ = &\quad \{\text{calculation above}\} \\ &\text{filter } (\text{all nd} \cdot \text{boxs}) \cdot \text{filter } (\text{all nd} \cdot \text{cols}) \cdot \text{filter } (\text{all nd} \cdot \text{rows}) \cdot \text{mcp} \cdot \\ &\quad \text{pruneBy rows} \\ = &\quad \{\text{interchanging the order of the filters}\} \\ &\text{filter } (\text{all nd} \cdot \text{rows}) \cdot \text{filter } (\text{all nd} \cdot \text{boxs}) \cdot \text{filter } (\text{all nd} \cdot \text{cols}) \cdot \text{mcp} \cdot \\ &\quad \text{pruneBy rows} \\ = &\quad \{\text{using the calculation above again}\} \\ &\text{filter } (\text{all nd} \cdot \text{rows}) \cdot \text{filter } (\text{all nd} \cdot \text{boxs}) \cdot \text{filter } (\text{all nd} \cdot \text{cols}) \cdot \text{mcp} \cdot \\ &\quad \text{pruneBy cols} \cdot \text{pruneBy rows} \\ = &\quad \{\text{repeating the last two steps one more time}\} \\ &\text{filter } (\text{all nd} \cdot \text{rows}) \cdot \text{filter } (\text{all nd} \cdot \text{boxs}) \cdot \text{filter } (\text{all nd} \cdot \text{cols}) \cdot \text{mcp} \cdot \\ &\quad \text{pruneBy boxs} \cdot \text{pruneBy cols} \cdot \text{pruneBy rows} \\ = &\quad \{\text{definition of } \text{filter correct}\} \\ &\text{filter correct} \cdot \text{mcp} \cdot \text{pruneBy boxs} \cdot \text{pruneBy cols} \cdot \text{pruneBy rows} \end{aligned}$$

Hence, we can define *prune* by

$$\begin{aligned} \text{prune} &:: \text{MatrixChoices} \rightarrow \text{MatrixChoices} \\ \text{prune} &= \text{pruneBy boxs} \cdot \text{pruneBy cols} \cdot \text{pruneBy rows} \end{aligned}$$

Readers who gave this solution (or a similar one in which the three components appear in any other order) can award themselves full marks.

The revised definition of *sudoku* now reads

$$\begin{aligned} \text{sudoku} &:: \text{Board} \rightarrow [\text{Board}] \\ \text{sudoku} &= \text{filter correct} \cdot \text{mcp} \cdot \text{prune} \cdot \text{choices} \end{aligned}$$

However, this version remains non-executable in practice. Again, assuming about a half of the 81 entries are fixed and an average of 3 choices/cell is generated by refining choices, there are still 3^{40} , or 12157665459056928801 boards to check. We still need something better.

4 One choice at a time

Humans employ a number of devices for filling in entries when solving Sudoku problems. For example, after pruning a matrix of choices, one or more entries that were previously blank may become fixed. In such a case, we can always prune again to see if more entries are filled in. The calculation above shows that we can have the composition of as many *prune* functions as we like. This is the way the simplest puzzles are solved.

There are also other strategies. For example, again after pruning the choice matrix it may turn out that a single row (or column or box) contains, for example, three entries such as 12, 12 and 123. It is clear that the third entry has to receive 3; if it receives 1 or 2, the first two entries cannot be filled in.

Repeatedly pruning the choice matrix is sensible, but we can combine it with another basic strategy. Rather than applying *mcp* when pruning fails to produce anything new, we can focus on one cell that has at least two choices, and generate a list of matrices in which this cell alone is expanded to each of its possible fixed choices.

Suppose we define a function

$$\text{expand} :: \text{Matrix Choices} \rightarrow [\text{Matrix Choices}]$$

that installs the fixed choices for one cell. This function satisfies the property that

$$\text{mcp} \approx \text{concat} \cdot \text{map mcp} \cdot \text{expand}$$

where \approx means equality up to a permutation of the answer. After applying *prune*, we can apply *expand* and then apply *prune* again to each of the results. Provided we discard any matrix that becomes blocked (see below), this process can be continued until we are left with a list of matrices, all of whose choices are fixed choices.

4.1 Blocked matrices

A matrix of choices can be **blocked** in that:

- One or more cells may contain zero choices. In such a case *mcp* will return an empty list;
- The same fixed choice may occur in two or more positions in the same row, column or box. In such a case *mcp* will still compute all the completed boards, but the correctness test will throw all of them away. [nodups on fixed fails](#)

Blocked matrices can never lead to a solution. In following the strategy of repeatedly pruning and expanding the matrix of choices, we can identify and discard any blocked matrix. Provided we do this, any remaining matrix that consists solely of fixed choices will be a solution to the puzzle.

Formally, we define *blocked* by

$$\begin{aligned}
 \text{blocked} &:: \text{Matrix Choices} \rightarrow \text{Bool} \\
 \text{blocked } cm &= \text{void } cm \vee \text{not } (\text{safe } cm) \\
 \text{void} &:: \text{Matrix Choices} \rightarrow \text{Bool} \\
 \text{void} &= \text{any } (\text{any null}) \\
 \text{safe} &:: \text{Matrix Choices} \rightarrow \text{Bool} \\
 \text{safe } cm &= \text{all } (\text{nodups} \cdot \text{fixed}) (\text{rows } cm) \wedge \\
 &\quad \text{all } (\text{nodups} \cdot \text{fixed}) (\text{cols } cm) \wedge \\
 &\quad \text{all } (\text{nodups} \cdot \text{fixed}) (\text{boxs } cm)
 \end{aligned}$$

4.2 Smallest number of choices

A good choice of cell on which to perform expansion is one with the *smallest* number of choices (greater than one of course). We will need a function that breaks up a matrix on the first entry with the smallest number of choices. A matrix that is not blocked is broken into five pieces:

$$cm = \text{rows1} \uparrow\uparrow [\text{row1} \uparrow\uparrow cs : \text{row2}] \uparrow\uparrow \text{rows2}$$

The smallest-choice entry is *cs*. The definition of *expand* is

$$\begin{aligned}
 \text{expand } cm &= [\text{rows1} \uparrow\uparrow [\text{row1} \uparrow\uparrow [c] : \text{row2}] \uparrow\uparrow \text{rows2}] \mid c \leftarrow cs \\
 \text{where } (\text{rows1}, \text{row} : \text{rows2}) &= \text{break } (\text{any best}) cm \\
 (\text{row1}, cs : \text{row2}) &= \text{break best row} \\
 \text{best } cs &= (\text{length } cs = n) \\
 n &= \text{minchoice } cm
 \end{aligned}$$

The definition of *minchoice* is

$$\text{minchoice} = \text{minimum} \cdot \text{filter } (> 1) \cdot \text{concat} \cdot \text{map } (\text{map length})$$

The number of choices in each cell is computed twice in the above definition, and it may be more efficient to avoid this duplication of effort. Also, we could probably make *minchoice* more efficient since once we have found an entry with two choices there is no point in looking further. Observe that *expand* returns \perp if there is no entry with at least two choices, for then *n* is undefined.

With this definition of *expand* we have

$$mcp \approx \text{concat} \cdot \text{map } mcp \cdot \text{expand}$$

Hence

$$\begin{aligned}
 &\text{filter correct} \cdot mcp \\
 \approx &\quad \{\text{above law of } \text{expand}\} \\
 &\text{filter correct} \cdot \text{concat} \cdot \text{map } mcp \cdot \text{expand} \\
 = &\quad \{\text{since } \text{filter } p \cdot \text{concat} = \text{concat} \cdot \text{map } (\text{filter } p)\}
 \end{aligned}$$

$$\begin{aligned}
& \text{concat} \cdot \text{map} (\text{filter correct} \cdot \text{mcp}) \cdot \text{expand} \\
= & \quad \{\text{property of } \text{prune}\} \\
& \text{concat} \cdot \text{map} (\text{filter correct} \cdot \text{mcp} \cdot \text{prune}) \cdot \text{expand}
\end{aligned}$$

Writing $\text{search} = \text{filter correct} \cdot \text{mcp}$ we therefore have

$$\text{search} = \text{concat} \cdot \text{map} (\text{search} \cdot \text{prune}) \cdot \text{expand}$$

This equation can be used as the basis of a definition of search provided we trap the terminal cases:

$$\begin{aligned}
& \text{search} :: \text{Matrix Choices} \rightarrow [\text{Matrix Choices}] \\
& \text{search } cm \\
& \quad | \text{ blocked } cm \quad = [] \\
& \quad | \text{ all (all single) } cm \quad = [cm] \\
& \quad | \text{ otherwise} \quad = (\text{concat} \cdot \text{map} (\text{search} \cdot \text{prune}) \cdot \text{expand}) cm
\end{aligned}$$

4.3 Final version

Now we can write down our final definition of sudoku :

$$\begin{aligned}
& \text{sudoku} :: \text{Board} \rightarrow [\text{Board}] \\
& \text{sudoku} = \text{map} (\text{map head}) \cdot \text{search} \cdot \text{prune} \cdot \text{choices}
\end{aligned}$$

The function $\text{map} (\text{map head})$ converts a matrix of singleton choices into a board. The program is quite fast, rarely taking more than a second or two to solve a puzzle.

5 Conclusions

The Sudoku problem provides an ideal classroom example with which to illustrate manipulations of arrays as well as manipulation of programs. Indeed, the pearl is more or less a straightforward transcription of two lectures I gave to first-year undergraduates, omitting most of the calculations. There is a temptation to identify array elements by their cartesian coordinates, and to think of the rows, columns and boxes of an array in terms of operations on these coordinates. Instead, and this is the pedagogic value of the exercise, we have gone for *wholemeal* programming, identifying these structures as complete entities in themselves. There are other Sudoku solvers out there, but the present one certainly seems one of the clearest and simplest.