



## Sudoku Players' Forums

[FAQ](#)
[Search](#)
[Memberlist](#)
[Usergroups](#)  
[Profile](#)
[You have no new messages](#)
[Log out \[ denis\\_berthier \]](#)

### THE REAL DISTRIBUTION OF MINIMAL PUZZLES

Goto page [Previous](#) [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#) [Next](#)



[Sudoku Players' Forums Forum Index](#) -> [General/puzzle](#)

[View previous topic](#) :: [View next topic](#)

#### Author

#### Message

**coloin**

Posted: Tue Jul 14, 2009 6:04 pm Post subject:



Indeed, its both simpler and a lot more complicated.

Joined: 06 May 2005  
Posts: 1045  
Location: Devon UK

There are MANY MANY unavoidable sets.....2-digit,3-digit,4-digit,5-digit etc. Some are massive [61 clues I believe]. The observer really needs to confirm this !

There is certainly more than 90 2-digit and 3-digit combined. Although yes you could pick 30 clues which cover all these and you might indeed have coincidentally covered the rest of the 3+ unavoidable. It would be unusual not to be able to find a 31st clue which defines the puzzle.

Essentially what a bottom up generator is doing [like Mikes] is picking selectively 20 or so [unconstrained] clues. What you might have then is 2000 or more say possible solution grids. [all with different remaining unavoidable sets]

Now.....if by chance.....one of those solution grids has "only" 2 unavoidable sets not covered...then a valid puzzle can be made by adding 2 clues. If the unavoidable sets have an overlapping clue - then the puzzle can be made/completed with 1 more clue. This happens with every clue in every 21 clue puzzle. !

Of course if one of the solution grids had many more still uncovered unavoidable sets then if the puzzle has to have that solution many more clues will be required. Superfluous clues will make the puzzle non-minimal.

This explains why the bottom up method with an undefined grid has a lower number of clues.

Essentially what a fast solver is doing is solving the puzzle and confirming [indirectly] that all the unavoidable sets in the solution grid are covered.

That all the unavoidable sets can be covered with <20 clues is quite an achievement. The fact that they can be covered in 17 clues in some grids maybe is testament to the statistical enormity of grids/unavoidable sets/clue selection in puzzles. There has to be minimum number of clues and it just happens to be 17 in some grids. Most grids have a 19 in them somewhere !

To **denis** - I have noted that hard puzzles tend to have many [ ? all] of the clues uniquely covering large unavoidable sets !

C

[Back to top](#)



**denis\_berthier**

Posted: Wed Jul 15, 2009 2:27 am Post subject:



Joined: 19 Jun 2007  
Posts: 730  
Location: Paris, France

#### coloin wrote:

A top down generator , by removing a clue and checking for sol>1 is actually verifying that an unavoidable set is not covered.  
The problem for random "puzzles" is dependant on the clue numbers  
Above 39 - all non-minimal  
31-39 - nearly all non-minimal, some >1 sol.  
29,30 - mostly non-minimal, most >1sol  
below 29 - almost all >1sol  
below 17 - all >1 sol.  
So prospects are hopeful only for 29 and 30 clue minimal puzzles !

I can see much has been done while I was away yesterday.

Do you have a more precise estimation of the % of minimal puzzles (wrt to the valid ones) for each number of clues.

If we had this, we could have a more precise idea of the top-down generators bias.

BTW, I've edited the end of my answer to **eleven's** mini example. The first version seemed too negative.

After all, the main part remains: all the valid puzzles with the same number of clues have the same probability of being reached by a top-down generator.

An additional correction factor must be added to my  $P_{n+1}/P_n$  formula, which depends on the % of minimals, but, if it is almost hopeless to compute it precisely, it may be possible to estimate it experimentally.

Coloin, beware: the result on the same probabilities is valid only when you start from complete grids. If you start from partial grids, some paths won't appear.

Last edited by denis\_berthier on Wed Jul 15, 2009 2:31 am; edited 1 time in total

[Back to top](#)

[profile](#) [pm](#) [www](#)

**denis\_berthier**

Posted: Wed Jul 15, 2009 2:30 am Post subject:

[quote](#) [edit](#)

Joined: 19 Jun 2007  
Posts: 730  
Location: Paris, France

**coloin wrote:**

To **denis** - I have noted that hard puzzles tend to have many [ ? all] of the clues uniquely covering large unavoidable sets !

Seems sensible. It's generally more difficult to draw conclusions from a large unavoidable set.

[Back to top](#)

[profile](#) [pm](#) [www](#)

**denis\_berthier**

Posted: Wed Jul 15, 2009 5:43 am Post subject:

[quote](#) [edit](#)

Joined: 19 Jun 2007  
Posts: 730  
Location: Paris, France

### **A TOP-DOWN GENERATOR WITH CONTROLLED BIAS**

I started this thread with the goal of:

- determining the bias of top-down generators;
- dealing with it in a non-standard way: instead of trying to modify the generator so that it has no bias (which seems an unrealistic goal), my idea was to introduce correction factors in the formulæ allowing to compute mean values.

This strategy led me to prove that all the valid puzzles with the same number of clues have the same probability of being reached by a top-down generator. This is more or less obvious, as the relative probabilities of 2 such puzzles depend only on the numbers of paths leading to each of them, which in turn don't depend on what is below any of the minimal puzzles at upper floors.

I also wrote a very simple formula:  $P_{n+1}/P_n = (n+1)/(81-n)$  that was supposed to give the successive ratios of these probabilities.

Due to an example by **eleven**, I realised that this formula neglected some "probability leakage" due to the multi-solution puzzles Q issued from a valid non minimal puzzle P, such that Q can also be considered as issued from a minimal puzzle at the same floor as P. A second correction factor should therefore be added. This factor is probably very close to 1 because the proportion of minimal puzzles is small at each floor, but there appears to be no easy means of computing it exactly.

I'll now propose **a modified top-down generator for which a simple formula for  $P_{n+1}/P_n$  holds exactly**, which will produce puzzles with more clues in the mean.

Consider first the classical top-down algorithm (for 1 puzzle):

**Code:**

```
1) generate a random complete grid
2) loop:
   let P be the current puzzle
   2a) choose one clue randomly from P and delete it, you get a puzzle P2
   2b) if P2 is minimal, return P2
   2c) if P2 has several solutions, GOTO 2a
```

```

2d) otherwise, set P=P2
end loop

```

Clause 2c is the cause of our probability leakage. Moreover, it is easy to see that it causes the generator to go deeper, i.e. towards puzzles with fewer clues.

Consider therefore the following modified algorithm (still for 1 puzzle):

**Code:**

```

1) generate a random complete grid
2) loop:
   let P be the current puzzle
   2a) choose one clue randomly from P and delete it, you get a puzzle P2
   2b) if P2 is minimal, return P2
   2c) if P2 has several solutions, GOTO 1
   2d) otherwise, set P=P2
end loop

```

The only difference is in case 2c: if we find a multi-sol puzzle, instead of backtracking to the previous state, we merely discard the current complete grid and restart the search with another one.

Notice that, contrary to the standard algorithm which produces a puzzle per complete grid, the modified algorithm will generally use several (the question is, how many?) complete grids before it outputs a puzzle.

I think someone proficient in C (**eleven??**) could easily do the modifications in the top-down version `suexgx.x` of `suexg` I used to generate the `sudogen0_1M` collection ([http://www.carva.org/denis.berthier/HLS/Classification/sudoku\\_gen.c](http://www.carva.org/denis.berthier/HLS/Classification/sudoku_gen.c)). Notice that the "official" `suexg14` wouldn't fit the purposes described above, as it is bottom up.

Maybe **Allan Barker** could also try this in his top-down generator.

The modified algorithm is likely to be much slower than the standard one, but it will be much faster than any unbiased generator (at least, given our current ways of imagining them) and the purpose here is not speed. To be clear: I don't mind having to let the program run for several days or weeks to get 10,000 minimal puzzles if I'm sure I can obtain correct statistics from them: this computation has to be done only once.

My interest in puzzle generators is recent (and limited to their usage for statistics). I'm not claiming the modified algorithm is new; I don't know. Perhaps, someone had the same idea previously. But, if such is the case, we now have good motivations for a new exploration of the possibilities.

Last edited by denis\_berthier on Fri Jul 17, 2009 9:38 am; edited 1 time in total

[Back to top](#)

 [profile](#)  [pm](#)  [www](#)

**eleven**

▢ Posted: Wed Jul 15, 2009 10:30 am Post subject:

 [quote](#)

**denis\_berthier wrote:**

Coloin, beware: the result on the same probabilities is valid only when you start from complete grids. If you start from partial grids, some paths won't appear.

Yes, of course.

From different 40 clue subgrids containing the same two 20 clues you will get different ratio's to find them. As i already had noticed, the probability to get to *non minimal* puzzles with n clues from top is not the same for all of the same level.

But i still believe, its the same for minimal puzzles. Then there should not be a bias between the puzzles of the same level and - given we know the correct clue distribution - the top-down can be used to create unbiased puzzle collections.

**denis\_berthier wrote:**

I think someone proficient in C (**eleven??**) could easily do the modifications in the top-down version `suexgx.x` of `suexg` I used to generate the `sudogen0_1M` collection

This change should do it (m0: is, where a new grid is created)

**Code:**

```

/*
for (i1=1;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())>1)A[P[i1]]=s1;}}
*/
for (i1=1;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())>1)goto m0;}}

```

But after counting for each level, how much multi solution puzzles are met on the way, i guess that this change will make the generation 15mio times slower.

**Edit:** oops, dont use that, it cant arrive at a solution.

Last edited by eleven on Wed Jul 15, 2009 1:28 pm; edited 1 time in total

[Back to top](#)

[profile](#) [pm](#)

**ronk**

Posted: Wed Jul 15, 2009 11:04 am Post subject:

[quote](#)

**denis\_berthier wrote:**

Coloin, beware: the result on the same probabilities is valid only when you start from complete grids. If you start from partial grids, some paths won't appear.

If, for  $n \geq 40$ , all  $n$ -clue subgrids are non-minimal, how can starting with a 40-clue subgrid lead to misleading results ?

Last edited by ronk on Wed Jul 15, 2009 11:06 am; edited 1 time in total

[Back to top](#)

[profile](#) [pm](#)

**denis\_berthier**

Posted: Wed Jul 15, 2009 11:06 am Post subject:

[quote](#) [edit](#)

**eleven wrote:**

the probability to get to *non minimal* puzzles with  $n$  clues from top is not the same for all of the same level.

The probability of all the *valid* (minimal or not) puzzles with  $n$  clues *is* the same - provided that you start from all the complete grids. Your mini-example led me to realise that a correction was necessary in my  $P_{n+1}/P_n$  formula, because of the probability leakage. But this correction is the same for all the valid puzzles at floor  $n$ . You can check this by retracing all the steps of my answer to your mini-example.

**eleven wrote:**

But i still believe, its the same for minimal puzzles. Then there should not be a bias between the puzzles of the same level and - given we know the correct clue distribution - the top-down can be used to create unbiased puzzle collections.

But

- 1) do we know the correct clue distribution?
- 2) even if we do, how can we use it to get unbiased collections? I think any unbiased algorithm will be much slower than my modified version of top-down suexg.

**eleven wrote:**

This change should do it (m0: is, where a new grid is created)

**Code:**

```
/*
for (i1=1; i1<=81; i1++) {s1=A[P[i1]]; if (s1) {A[P[i1]]=0; if (solve())>1 A[P[i1]]=s1;}}
*/
for (i1=1; i1<=81; i1++) {s1=A[P[i1]]; if (s1) {A[P[i1]]=0; if (solve())>1 goto m0;}}
```

But after counting for each level, how much multi solution puzzles are met on the way, i guess that this change will make the generation 15mio times slower.

Thank you. I changed this line and compiled it.

I think you're right about the increased complexity.

I've launched it 10 mn ago, but haven't yet got any puzzle. I'll let it run as long as possible.

[Back to top](#)

[profile](#) [pm](#) [www](#)

**denis\_berthier**

Posted: Wed Jul 15, 2009 11:08 am Post subject:

[quote](#) [edit](#)

**ronk wrote:**

Joined: 19 Jun 2007

Joined: 19 Jun 2007

Posts: 730

Location: Paris, France

**denis\_berthier wrote:**

Coloin, beware: the result on the same probabilities is valid only when you start from complete grids. If you start from partial grids, some paths won't appear.

If, for  $n \geq 40$ , all  $n$ -clue subgrids are non-minimal, how can starting with a 40-clue subgrid lead to misleading results ?

My sentence was ambiguous: in my forest or in the case of a top-down generator, you start (at least virtually) from ALL the complete grids.

The result is true also for all the valid puzzles issued from a single complete grid.

But, if you consider only the descendants of a 40-clue valid subgrid, for some of them paths that would have come from other parts of the complete grid by deleting clues in a different order will be missing. These puzzles will miss part of their heritage.

[Back to top](#)
[profile](#) [pm](#) [www](#)
**Allan Barker**

Posted: Wed Jul 15, 2009 11:58 am Post subject:

[quote](#)

Joined: 21 Feb 2008

Posts: 290

Location: Bangkok

**denis\_berthier wrote:**

Maybe Allan Barker could also try this in his top-down generator.

That's easy to do but the results are not promising. Results from 100000 trials showed a single puzzle, two more runs produced nothing. I normally make 100 to 300 puzzles/sec. so the estimated time to make 10000 puzzles > 100 days.

**Code:**

```

30      0
29      0
28      0
27      0
26      0
25      1
24      0
23      0
22      0
21      0
20      0
19      0
18      0
-----
N=1 avg=25.00

```

Maybe a bit more interesting is the following data where a top down generator records all the puzzles it encountered on the way down, recorded per clue size. Clearly it starts running into lots of multi-solution puzzles long before minimals are possible.

**Code:**

```

Puzzles seen on the way down
  minimal    multi non-min    (minimal/total)
56          0         0         0  0.0000
55          0         0  100000  0.0000
54          0         69  100000  0.0000
53          0        171  100000  0.0000
52          0        290  100000  0.0000
51          0        487  100000  0.0000
50          0        654  100000  0.0000
49          0        939  100000  0.0000
48          0       1338  100000  0.0000
47          0       1705  100000  0.0000
46          0       2234  100000  0.0000
45          0       2909  100000  0.0000
44          0       3585  100000  0.0000
43          0       4839  100000  0.0000
42          0       6219  100000  0.0000
41          0       7995  100000  0.0000
40          0       9618  100000  0.0000
39          0      12256  100000  0.0000

```

38	0	15068	100000	0.0000
37	0	19252	100000	0.0000
36	0	23892	100000	0.0000
35	0	29468	100000	0.0000
34	0	37269	100000	0.0000
33	0	46616	100000	0.0000
32	0	58325	100000	0.0000
31	0	74271	100000	0.0000
30	0	96798	100000	0.0000
29	13	124509	99987	0.0001
28	271	165236	99716	0.0010
27	2451	225735	97265	0.0075
26	11872	310259	85393	0.0291
25	29689	396323	55704	0.0616
24	34407	402155	21297	0.0751
23	17573	257384	3724	0.0631
22	3468	85695	256	0.0388
21	250	12452	6	0.0197
20	6	622	0	0.0096
19	0	15	0	0.0000
18	0	0	0	0.0000
-----				
N=100000 24.37				

[Back to top](#)**denis\_berthier**

Posted: Wed Jul 15, 2009 12:02 pm Post subject:



Joined: 19 Jun 2007  
 Posts: 730  
 Location: Paris, France

**Allan Barker wrote:****denis\_berthier wrote:**

Maybe Allan Barker could also try this in his top-down generator.

That's easy to do but the results are not promising. Results from 100000 trials showed a single puzzle, two more runs produced nothing. I normally make 100 to 300 puzzles/sec. so the estimated time to make 10000 puzzles > 100 days.

That's much faster than suexg and very interesting.  
 Is your generator available?  
 I could let it run for a a long time.

[Back to top](#)**eleven**

Posted: Wed Jul 15, 2009 12:28 pm Post subject:



Joined: 10 Feb 2008  
 Posts: 358

**eleven wrote:**

This change should do it (m0: is, where a new grid is created)

**Code:**

```
/*
for (i1=1;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())>1A[P[i1]]=s1;}}
*/
for (i1=1;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())>1goto m0;}}
```

Sorry, this change is rubbish and my time estimation either.  
 sudogen calculates a random order, in which it tries to eliminate the clues, P[i].  
 With this change (when getting multi solutions, restart), it for sure will restart each time.  
 Maybe later i will have the time to correct it.

[Back to top](#)**Pat**

Posted: Wed Jul 15, 2009 12:57 pm Post subject:



Joined: 18 Jul 2005  
 Posts: 1471

**coloin wrote:**

There are MANY MANY unavoidable sets.....2-digit,3-digit,4-digit,5-digit etc.  
 Some are massive [61 clues I believe].

the largest i've seen is **60** clues,  
 found by **Ocean** (2006.Nov.23)

~ Pat

[Back to top](#)

coloin

Posted: Wed Jul 15, 2009 3:29 pm Post subject:

**@Pat** - yes its 60...btw - no one has said "\*\*\*\* arnt there an awful lot of unavoidable yet 🙄"<http://www.sudoku.com/boards/viewtopic.php?t=4771&postdays=0&postorder=asc&start=30>**JPF** got 50 % >1 sol at around the 43 clue mark - same as **Allan**I cant really see minimal puzzles from **Ocean** data.But congratulations to **Allan** on the first random puzzle ever !!**@denis** - I was considering the generation of 40-clue subgrids - with no superfluous clues - surely this doesnt suffer from losses ?.....this actually means there are **2** clues in each unavoidable set.....

C

[Back to top](#)

David P Bird

Posted: Wed Jul 15, 2009 4:35 pm Post subject:



If our target solution grid contains a UR pattern

a b

b a

we know that any set of random exclusions that eliminates all four of these cells will produce a puzzle with alternative solutions. The unavoidable set checking system I described yesterday allows random selections to be made until it is detected that only one of these cells still exists in the reduced puzzle. Which one of the four survives is left to chance, but at that stage we know that any further random choice that lands on that cell will lead to failure, so we take it out of the mix by protecting it.

Now working out the probabilities of any cell being protected this way at each of the random steps is horrendous, which defeats the purpose of this thread. However, if we ran the algorithm enough times we could build up average values for the number of protected cells at each step. Whether or not this is worth doing depends on how well these heuristic values could be woven into a probability calculation. At this point I gracefully exit stage left!

[Back to top](#)

eleven

Posted: Wed Jul 15, 2009 4:36 pm Post subject:



Denis, this should work now, i tried to keep dukuso's style 😊

**Code:**

```
Replace the line
for(i1=1;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())>1A[P[i1]]=s1;}}
by
for(i1=1;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())>1){A[P[i1]]=s1;break;}}
i=++i1;for(i1=i;i1<=81;i1++){s1=A[P[i1]];if(s1){A[P[i1]]=0;if(solve())<2)goto m0;
A[P[i1]]=s1;}}
```

Note that there are big speed differences between PC's. On an "Intel(R) Core(TM)2 Duo CPU E8400 @ 3.00GHz" running 2 instances of sudo\_gen (with different seed) i can generate about 180 puzzles each per second. For the new version it took about 10-15 minutes to get a puzzle (i got 2 26's, a 27 and 2 28's). So with full speed on such a PC you could get ~10/hour. But if you only get 20 puzzles/sec, you would need 18 times longer, not much more than 10/day.

Maybe Allan can run sudo\_gen on his machine to compare the speed.

[Back to top](#)Display posts from previous:   [Sudoku Players' Forums Forum Index -> General/puzzle](#)All times are GMT  
Goto page [Previous](#) [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#) [Next](#)

[Stop watching this topic](#)

Jump to:

You **can** post new topics in this forum  
You **can** reply to topics in this forum  
You **can** edit your posts in this forum  
You **can** delete your posts in this forum  
You **can** vote in polls in this forum

Powered by phpBB © 2001, 2005 phpBB Group