



A Pseudo-Random Number Generator for Spreadsheets

Michael Lampton, Space Sciences Lab, UC Berkeley

Abstract

Setting up a spreadsheet to simulate noisy data collection from an experiment requires a generator of pseudo-random numbers. The function `RAND` built into popular spreadsheets is unsuitable because it rerandomizes every time the spreadsheet is recalculated: unlike measured data, `RAND` changes for each data analysis activity. Moreover, spreadsheet rootfinding plugins necessarily recalculate the entire spreadsheet for each internal iteration, and `RAND` cannot keep its output constant during an iterative fit to the simulated data. Beyond that there is no way to verify agreement between instructor and student statistics. Here, I introduce a well-tested random number generator that overcomes these limitations and I show how to make it portable to spreadsheets and high-level computer languages.

Keywords: random numbers, spreadsheets.

1. Goals

For decades, physics instructors have used spreadsheets to organize computational work in the classroom and laboratory. To statistically model measurement errors, some kind of pseudo-random number (PRN) generator is essential. Desirable features of a spreadsheet-based PRN generator (PRNG) are:

- Occupies a single spreadsheet cell;
- Specifies a single input cell for its argument or seed;
- Output spans $0 < x < 1$ for use with distribution generators `NORMINV`, `TINV`, etc.;
- Has a period that far exceeds the likely number of rows in an experiment simulator;
- Has correct means, variances, uniformities, and autocorrelations;

- Has a white (flat) power spectrum;
- Pseudo-random sequences (PRSs) from different seeds should be statistically orthogonal;
- For the same seed, different spreadsheets should yield the same PRS;
- Should also be portable between spreadsheets and high level languages.

2. A fully portable generator

The most popular PRN generators are based on the multiplicative linear algorithm originally by [Lehmer \(1949\)](#); see also [Wichmann and Hill \(1982\)](#), [Park and Miller \(1988\)](#), [Press, Flannery, Teukolsky, and Vetterling \(1988\)](#), [Knuth \(1997\)](#), [L'Ecuyer \(1998\)](#), and [Wichmann and Hill \(2006\)](#). This generator receives an integer Zin and produces an integer $Zout$:

$$Zout = (A \cdot Zin) \bmod M \quad (1)$$

Briefly, the low order bits of the input integer Zin are hoisted to higher significance by an integer multiplier A , and the highest order bits are discarded by the modulus operation. The modulus M sets the finesse of the comb of output values. If M and A are properly chosen, the period of the sequence will have the maximum possible length of $M - 1$. Owing to its speed and simplicity it has enjoyed a long life and has been ported to a variety of environments. Obviously it must never be given a seed of zero, or the whole sequence would collapse. Indeed, seed values very near zero cause the first few iterates to be substandard in size (see [Wichmann and Hill \(2006\)](#)), and seeds must be randomized if the PRS is to have a random-looking startup sequence. I offer a seeder in the next section.

Exact integer arithmetic is essential if a generator is to be portable. The IEEE-754 double precision specification requires exact integer arithmetic in the range -2^{52} to $+2^{52} - 1$ and this standard is widely obeyed by high-level languages. However, spreadsheets typically deliver exact integer arithmetic only for numbers whose size is less than about 10^{15} , i.e., fifteen digits accuracy or about 2^{50} (see [Almiron *et al.* \(2010\)](#)). To avoid overflow, a PRN algorithm internal integer product $A \cdot M$ should be comfortably smaller than this fifteen-digit bound.

A second requirement is that the generator must deliver variates in the range $0 < Z < 1$ so that these can be fed into the appropriate inverse distribution generator, for example `NORMINV` or `TINV` spreadsheet functions. This requirement is customarily met by moving Equation 1 into a floating point environment, with $M \cdot Zin$ playing the role of the integer Zin :

$$Zout = ((M \cdot A \cdot Zin) \bmod M) / M \quad (2)$$

I adopt the “minimal standard generator” of [Park and Miller \(1988\)](#) defined by the constants

$$\begin{aligned} M &= 2^{31} - 1 = 2147483647 \\ A &= 7^5 = 16807 \end{aligned} \quad (3)$$

This generator has been exhaustively tested and has been found to be generally acceptable within its limited sequence length of $M - 1$. This generator will be portable among systems whose arithmetic correctly handles numbers of magnitude $A \cdot M$ without loss of integer accuracy. Here, the largest number that has to be handled is the MA product, about $3.6 \cdot 10^{13}$. This is comfortably within the 15-digit accuracy of all popular spreadsheets. Features of this PRNG are:

- Successive iterates are uncorrelated;
- Beyond that, they pass the Knuth spectral test for dimensions 2,3,4,5, and 6;
- The PRNG is maximal length (here equal to $M - 1$);
- The PRNG populates the variate axis uniformly;
- The PRNG is portable: it delivers the same PRS on every platform.

In a spreadsheet cell, Equation 2 is implemented with the expression

$$=\text{MOD}(\text{ROUND}(M*A*Z,0),M)/M \quad (4)$$

Again, M and A are the constants from Equation 3. Z stands for the address (column,row) of whichever cell contains the previous RN in your sequence, or, for the initial RN, it is the address of the seed generator. The $\text{ROUND}(X,0)$ operation is essential for portability: it reestablishes the correct integer product by removing the floating point division errors (order of 10^{-15}) that arise in the previous iteration division. Without $\text{ROUND}()$, the PRS would depend on the chain of division errors which would introduce fractional terms into the sum. These fractional terms differ among spreadsheets and can also differ from the IEEE-754 floating point specifications, causing the various platform PRSs to diverge after some iterations.

3. Portable seeding

How should this function be seeded? In experiment modeling, users will want a selection of seeds that deliver independent statistics. The entire PRS period is about two billion, so if each segment has ~ 1000 iterates, any randomly chosen seed will have only one chance in a million of overlapping any other given seed segment. Very good odds! but of course the RNG is totally deterministic and the seeder must be verified for freedom from overlap. Fractional seeds are essential since all integer inputs are equivalent to zero seed and yield the same null PRS.

A good seeder will accept an integer run number and deliver a prandomized seed value. (A simple list of seeds like 0.1, 0.2, etc., would fail the portability test because they are unlikely to be found among the comb of PRS values.) A prandomized seed makes the PRNG self-starting, requiring no warm-up iterations before use. I address this issue here by offering an explicit $\text{seed}()$ function that accepts an integer run number $\text{IRUN}=1, 2, \dots$ and delivers a starting point for a PRS. For portability, I precondition the seed by applying the same modulus treatment that each RN has:

$$=\text{MOD}(\text{ROUND}(\text{MOD}(\text{IRUN}*\text{EXP}(1),1)*M*A,0),M)/M \quad (5)$$

Here, the constant $\text{EXP}(1)=2.71828\dots$ supplies some fractional digits that are boosted by the integer run number IRUN . That fractional part is then boosted by the MA product to fill the working span of double precision integers. That product is then reduced modulo M , and normalized to unit span in the same way that PRS numbers are reduced, so that each seed is a member of the PRS comb. These actions make every seed compute the same way on all platforms.

4. Testing for portability

Park and Miller (1988) emphasize that exhaustive statistical testing is exceedingly demanding of resources (see for example Fishman and Moore (1986)) and recommend that any portable RNG should be tested for correctness rather than for its statistics. For the constants in Equation 3 above and a seed derived from IRUN=1, iteration 10000 should yield $Z = 0.785320384794$. I verified this result for Microsoft **Excel** 2007 (PC edition), Gnome **Gnumeric** 1.10.16, Open Office **Calc** 3.3.0, Java 1.5.0, Gnu C, and Python 2.6.6. Portions of these runs are listed in Table 1. All platforms tested are in agreement over the range of parameters tested, giving good evidence of portability.

IRUN=	1	2	3	999
seed=	0.162690911052	0.325381822570	0.488072733622	0.528220262159
iter=1	0.346142053300	0.692291932969	0.038433986268	0.797946102357
iter=2	0.609489807212	0.350517402566	0.960007209778	0.080142321568
iter=3	0.695189804628	0.145984931451	0.841174736079	0.951998594195
iter=4	0.055046384714	0.568742901352	0.623789286066	0.240372629482
Excel 10000	0.785320384794	0.056613301419	0.841933686213	0.887922685076
Gnumeric 10000	0.785320384794	0.056613301419	0.841933686213	0.887922685076
Calc 10000	0.785320384794	0.056613301419	0.841933686213	0.887922685076
Java 10000	0.785320384794	0.056613301419	0.841933686213	0.887922685076

Table 1: Demonstration of portability among platforms. For run numbers listed, PRS iterations 1, 2, 3, 4, and 10000 are shown to 12 digits accuracy. Iteration 10000 is shown for **Excel** 2007 (PC); **Gnumeric** (PC); OpenOffice3 **Calc**(Mac); and Java which is IEEE-754 compliant.

5. Testing the seeder

Any PRS run should not overlap any portion of any other PRS run. The seeder generates a list of random-appearing seeds that start the PRSs. Figure 1 shows the log of the maximum nonoverlapping PRS length as a function of the log of the number of runs (blue) and also (red) the maximum PRS length if the segments were to be uniformly arranged, head-to-tail, without gaps.

6. Demonstrating the RNG statistics by histogram

The statistical uniformity and scatter of the RNG can be shown with a simple histogram showing the accumulated hits in a number of equal size bins that span the variate range $0 < Z < 1$. In Figure 2, I have set up 100 equal bins in Java, and overplotted 100 runs of 60000 iterations. Each run should average about 600 hits per bin. The accumulated statistics show good homogeneity and a Gaussian distribution of hit densities.

7. Demonstrating the RNG statistics by power spectrum

A Fourier analysis of the each PRS can in principle reveal subtle periodicities that would escape binning tests. In Figure 3, I show power spectra computed in Java of 4096-length PRS vectors for IRUN=1 and IRUN=2. The Fourier power values are expected to be distributed with

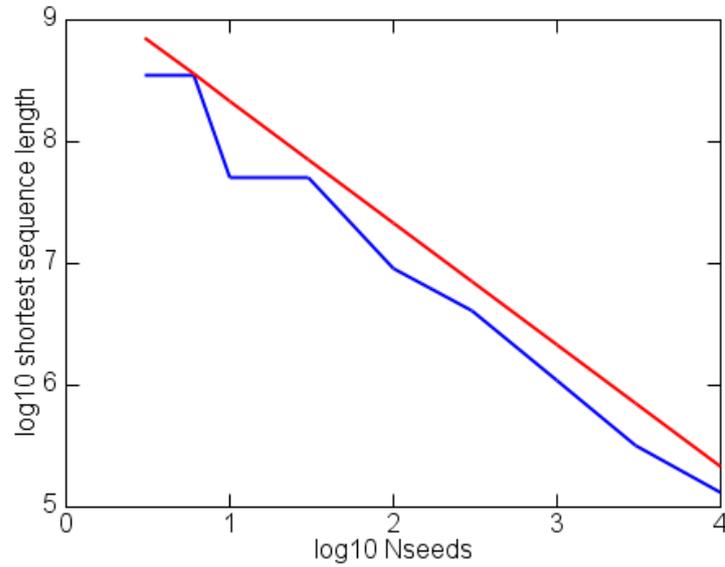


Figure 1: Plot showing the log of the maximum available sequence length without overlapping any other sequence as a function of the log of the number of seeds chosen. For 100 runs, individual sequences can have millions of iterations without overlap.

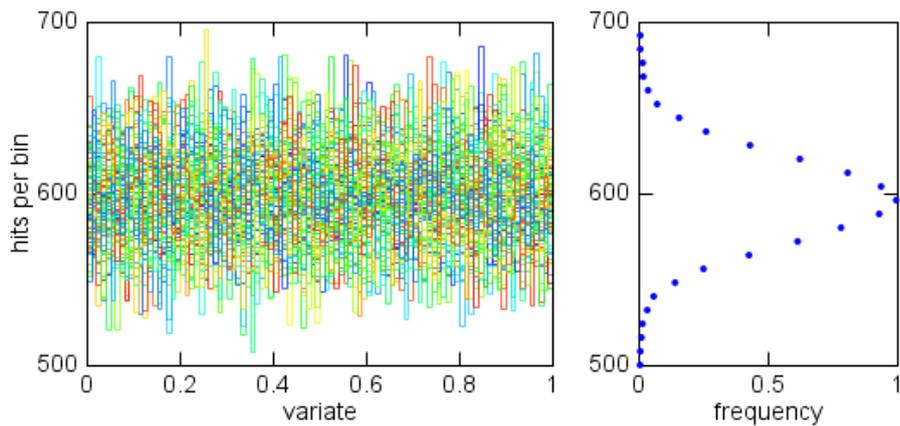


Figure 2: For a set of 100 equal bins spanning $0 < z < 1$, the number of RNs falling into each bin is plotted. 60000 iterations per run; 100 runs. Right: bin probabilities.

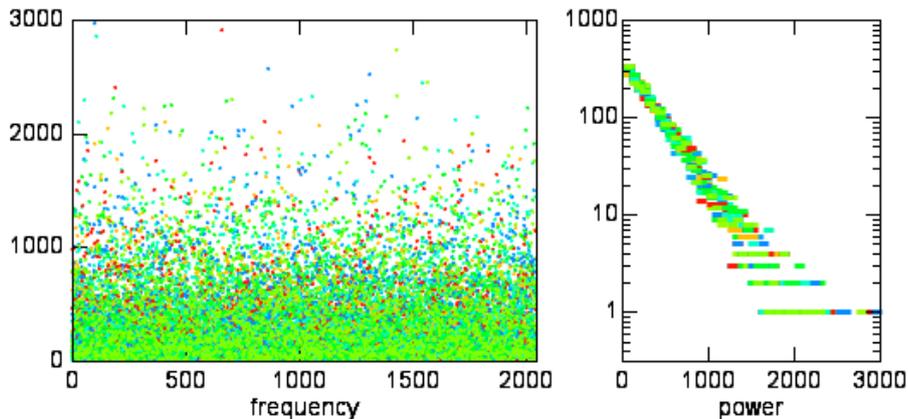


Figure 3: Left: Fourier power at frequencies 1 to 2048 cycles per sequence, for IRUN=1 (red) and IRUN=2 (blue). Right: histogram of power values. A negative exponential distribution would average to a straight descending line.

a negative exponential distribution. Figure 3 shows that this is the case, without strongly grouped or peaked frequency features and with the expected power distribution.

8. A handy on-off switch

I noted that IRUN=0 and its seed value of zero is prohibited because zero collapses the entire PRS. If however each PRN is used solely to feed an inverse cumulative probability function such as NORMINV, this all-zero off state can be recognized and employed to turn off the random deviations throughout the worksheet. Use the expression

$$=IF(Z=0,0,NORMINV(Z,0,1)) \quad (6)$$

where Z is again the (column, row) address of the uniform PRN being converted.

9. Conclusions

In Equations 4 and 5, I have presented a portable uniform random number generator that is under user control: its variates are effectively random and independent for a wide range of seed values, but then — like real measurements — they become constant during subsequent data analysis. Background theory and two simple demonstrations show good performance for small scale simulation situations found in the classroom and laboratory.

Acknowledgment

The author gratefully acknowledges the support by the Director, Office of Science, U.S. Department of Energy under Contract No. DE-AC03-76SF00098.

References

- Almiron MG, *et al.* (2010). “On the Numerical Accuracy of Spreadsheets.” *J. Stat. Software*, **34**(4), 1–24.
- Fishman GS, Moore LR (1986). “An Exhaustive Analysis of Multiplicative Congruential Random Number Generators with Modulus $2^{31} - 1$.” *SIAM J. Scientific and Statistical Computing*, **7**, 24–45.
- Knuth DE (1997). *The Art of Computer Programming*, volume 2: Seminumerical Algorithms. 3rd edition. Addison Wesley, Reading, Massachusetts.
- L’Ecuyer P (1998). *Handbook on Simulation*, chapter 4. Random Number Generation. 3rd edition. Wiley.
- Lehmer DH (1949). “Mathematical Methods in Large-scale Computing Units.” *Proc. 2nd Symposium on Large-Scale Digital Calculating Machinery*, pp. 141–146.
- Park SK, Miller KW (1988). “Random Number Generators: Good Ones are Hard to Find.” *Communications of the ACM*, pp. 1192–1201.
- Press WH, Flannery BP, Teukolsky SA, Vetterling WT (1988). *Numerical Recipes*. 2nd edition. Cambridge University Press.
- Wichmann BA, Hill ID (1982). “Algorithm AS183: An Efficient and Portable Pseudo-Random Number Generator.” *Applied Statistics*, **31**(2), 188–190.
- Wichmann BA, Hill ID (2006). “Generating Good Pseudo-Random Numbers.” *Computational Statistics and Data Analysis*, **51**(3), 1614–1622.

Affiliation:

Michael Lampton
Space Sciences Lab
University of California
Berkeley CA 94720 USA
E-mail: mlampton@SSL.berkeley.edu
URL: www.mikelampton.com