

TME – Assignment: testing random number generators

The "random number generators" we find in all programming languages and libraries are in fact "pseudo-random number generators", based on one or another algorithm. These generators can be tested to check the quality of the sequences they provide.

In this exercise you will implement such a test.

- Read chapter 7 of "Numerical recipes":

- https://www.cec.uchile.cl/cinetica/pcordero/MC_libros/NumericalRecipesinC.pdf
- https://websites.pmc.ucsc.edu/~fnimmo/eart290c_17/NumericalRecipesinF77.pdf

- Choose a random number generator to work with: look for the details of its implementation (algorithm, etc.) and then implement a "bad one" (as discussed in the book).

- Select also a "good" random number generator from your programming language. Find out from the language documentation how this generator has been implemented and document it.

- Implement four tests of the generators to compare and validate them. The following references are suggested to select the tests:

- <https://www.random.org/analysis/>
- <https://csrc.nist.gov/Projects/Random-Bit-Generation/Documentation-and-Software/Guide-to-the-Statistical-Tests>
- <https://nvlpubs.nist.gov/nistpubs/legacy/sp/nistspecialpublication800-22r1a.pdf>
- <https://csrc.nist.gov/projects/random-bit-generation/documentation-and-software>
- <http://webhome.phy.duke.edu/~rgb/General/dieharder.php>
- <http://simul.iro.umontreal.ca/testu01/tu01.html>

- Justify the choice and then discuss the results and the comparison of the "good" and "bad" generators.

- Can you imagine a situation in a simulation/experiment where the faults of the "bad" generator could affect the outcome?

Deliver in this page a report describing the above activities; include any program you have used. Pack all the deliverables in a single zip file.

1 Creating random generators:

I've chosen to try 2 different pseudo-random number generators from Chapter 7 of Numerical recipes [1].

First I'll try the one from page 276, which is a straight forward linear congruential generator with $a = 1103515245$, $c = 12345$ and $m = 2^{32}$, which are not particularly good choices, but neither gross embarrassments (period $\approx 2^{30} \approx 10^9$):

```

1 #include <string.h>
2 #include <stdlib.h>
3 #include <stdio.h>
4 #include <math.h>
5
6
7 unsigned long next=1;
8
9 int ran0(void){
10     next = next*1103515245 + 12345;
11     return (unsigned int)(next/65536) % 32768;
12 }
13
14 void sran0(unsigned int seed){
15     next=seed;
16 }
17
18
19 int main(){
20
21     sran0(8989743); //We set the seed here
22     int N =100000;
23     double ran[N];
24
25     for(int i = 0; i<N; i++){
26         ran[i]= static_cast <float> (ran0()) / static_cast <float> (RAND_MAX);
27     }
28
29     FILE* fichero;
30     fichero = fopen("Random_numbers1.txt", "w" );
31     for ( int i = 0; i < N; i++){
32         fprintf(fichero, "%f \n", ran[i]);
33     }
34
35     return 0;
36
37 }
```

then I'll try this same implementation but with very bad values, such as the ones IBM used in it's RANDU routine $a = 65539$, $c = 0$ and $m = 2^{31}$, which later should give bad results in the tests.

And finally I also want to try from page 279 a Minimal standard generator that applies Schrage's algorithm, with $a = 7^5 = 16807$, $m = 2^{31} - 1 = 2147483647$, $q = 127773$ and $r = 2836$, which should be better than the first one (period $\approx 2^{31} - 2 \approx 2.1 \cdot 10^9$) :

```

1 #include <string.h>
2 #include <stdlib.h>
3 #include <stdio.h>
4 #include <math.h>
5
6 #define IA 16807
7 #define IQ 127773
```

```

8 #define IM 2147483647
9 #define IR 2836
10 #define AM (1.0/IM)
11
12
13 unsigned long next=1;
14
15 void sran0(unsigned int seed){
16     next=seed;
17 }
18
19 float ran0(void){
20     long k;
21     float ans;
22
23     k=next/IQ;
24     next=IA*(next-k*IQ)-IR*k;
25     if (next < 0 ){
26         next += IM;
27     }
28     if (next>IM){
29         next -= IM;
30     }
31     ans=AM*next;
32     return ans;
33 }
34 }
35
36
37
38 int main(){
39
40     sran0(8989743); //We set the seed here
41     int N =100000;
42     double ran[N];
43
44     for(int i = 0; i<N; i++){
45         ran[i]= ran0();
46     }
47
48
49     FILE* fichero;
50     fichero = fopen("Random_numbers.txt", "w" );
51     for ( int i = 0; i < N; i++){
52         fprintf(fichero, "%f \n", ran[i]);
53     }
54
55     return 0;
56 }
57 }

```

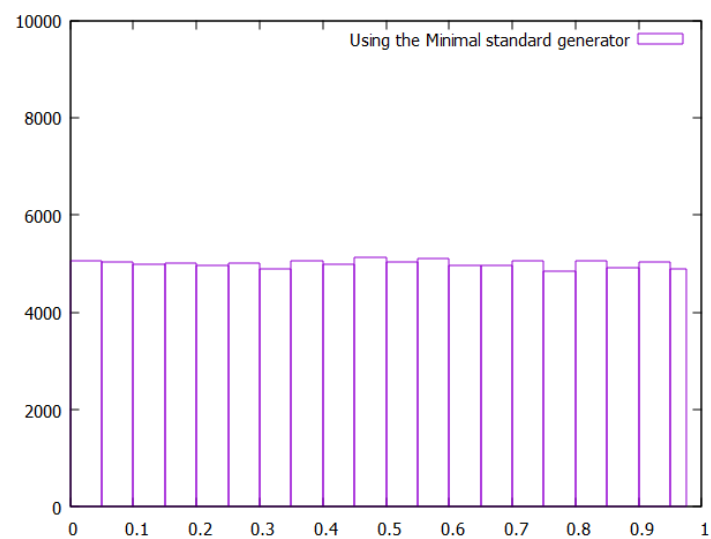
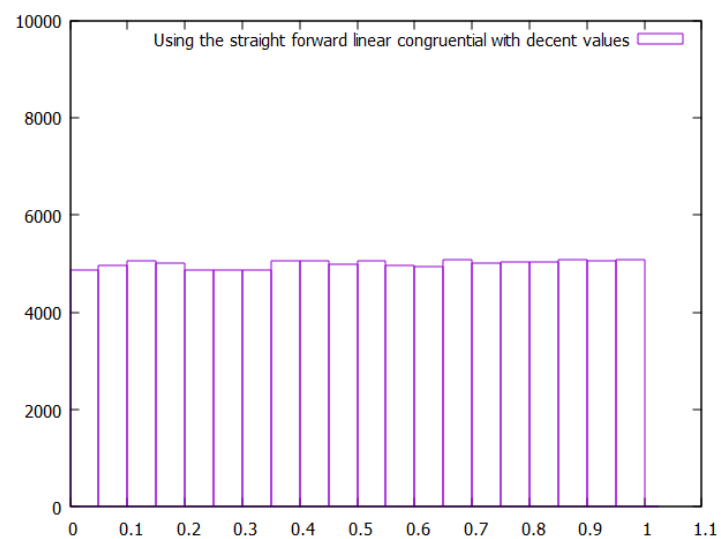
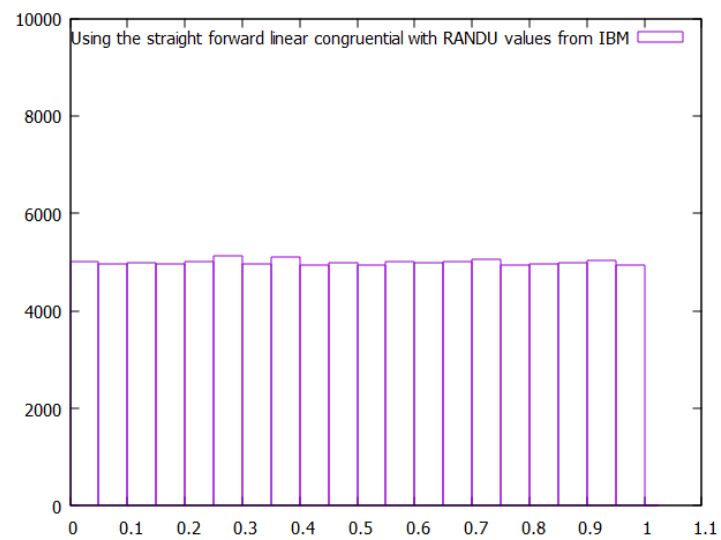
Which using this gnuplot code to plot the obtained data:

```

1 reset
2 cd 'C:\CARRERAS\MASTER UB\1r Semestre\Metodes Estadistics i matematics\
   Exercices\Random Numbers'
3
4 set yrange [0:100000]
5 binwidth=0.05
6 bin(x,width)=width*floor(x/width)
7
8 plot 'Random_numbers0.txt' using (bin($1,binwidth))+binwidth/2:(1.0) smooth
   freq with boxes

```

We get the following distributions for $N=100.000$ random numbers and a chosen seed = 8989743:



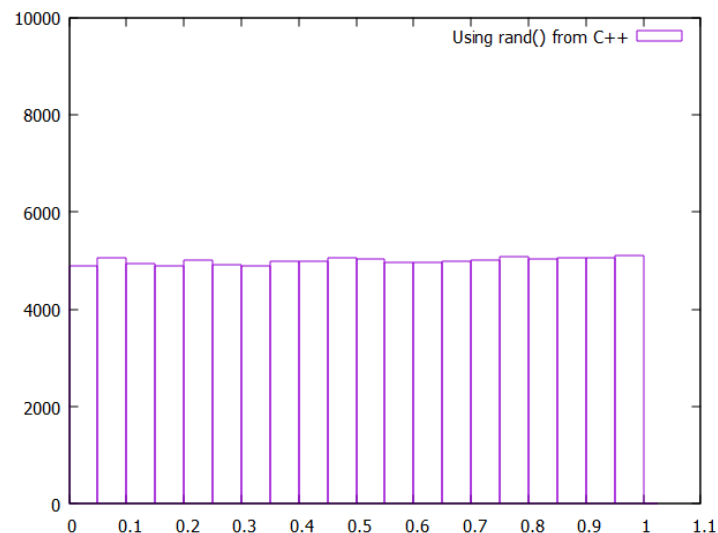
2 Choosing a "good" random number generator:

Now, we are going to choose the function `rand()`, imported from "`#include <stdlib.h>`", to do so, this time we don't need to define our random function, which is given, but instead we need to normalize it to give values from 0 to 1:

```

1 #include <string.h>
2 #include <stdlib.h>
3 #include <stdio.h>
4 #include <math.h>
5
6
7 int main(){
8
9     srand(8989743); //We set the seed here
10    int N = 100000;
11    double ran[N];
12
13    for(int i = 0; i < N; i++){
14        ran[i] = static_cast <float> (rand()) / static_cast <float> (RAND_MAX);
15    }
16
17    FILE* fichero;
18    fichero = fopen("Random_numbers2.txt", "w" );
19    for ( int i = 0; i < N; i++){
20        fprintf(fichero, "%f \n", ran[i]);
21    }
22
23    return 0;
24 }
```

Again, for the same values of N and the seed as before, this gives a distribution:



So, now we need to find its implementation, which after some research, we find that is given by:

```

1 /* Return a random integer between 0 and RAND_MAX. */
2 int
3 rand (void)
4 {
5     return (int) __random ();
6 }
```

Which calls the function `__random()`, that at the same time calls a function `__random()_r`:

```

1  /* If we are using the trivial TYPE_0 R.N.G., just do the old linear
2     congruential bit. Otherwise, we do our fancy trinomial stuff, which is
3     the same in all the other cases due to all the global variables that
4     have been set up. The basic operation is to add the number at the rear
5     pointer into the one at the front pointer. Then both pointers are
6     advanced to the next location cyclically in the table. The value
7     returned is the sum generated, reduced to 31 bits by throwing away the
8     "least random" low bit.
9
10    Note: The code takes advantage of the fact that both the front and rear
11    pointers can't wrap on the same call by not testing the rear pointer if
12    the front one has wrapped. Returns a 31-bit random number.  */
13
14  int
15  __random_r (buf, result)
16      struct random_data *buf;
17      int32_t *result;
18  {
19      int32_t *state;
20
21      if (buf == NULL || result == NULL)
22          goto fail;
23
24      state = buf->state;
25
26      if (buf->rand_type == TYPE_0)
27      {
28          int32_t val = state[0];
29          val = ((state[0] * 1103515245) + 12345) & 0x7fffffff;
30          state[0] = val;
31          *result = val;
32      }
33      else
34      {
35          int32_t *fptr = buf->fptr;
36          int32_t *rptr = buf->rptr;
37          int32_t *end_ptr = buf->end_ptr;
38          int32_t val;
39
40          val = *fptr += *rptr;
41          /* Chucking least random bit. */
42          *result = (val >> 1) & 0x7fffffff;
43          ++fptr;
44          if (fptr >= end_ptr)
45          {
46              fptr = state;
47              ++rptr;
48          }
49          else
50          {
51              ++rptr;
52              if (rptr >= end_ptr)
53                  rptr = state;
54          }
55          buf->fptr = fptr;
56          buf->rptr = rptr;
57      }
58      return 0;
59
60  fail:
61      __set_errno (EINVAL);
62      return -1;
63  }

```

Which at the end, we see is another linear congruential generator with $a = 1103515245$ and $c = 12345$ such as the first one we choose, that plays more with cycling pointers and the size of the values, to be more efficient and give a better pseudo-random generator.

3 Implement four tests of the generators

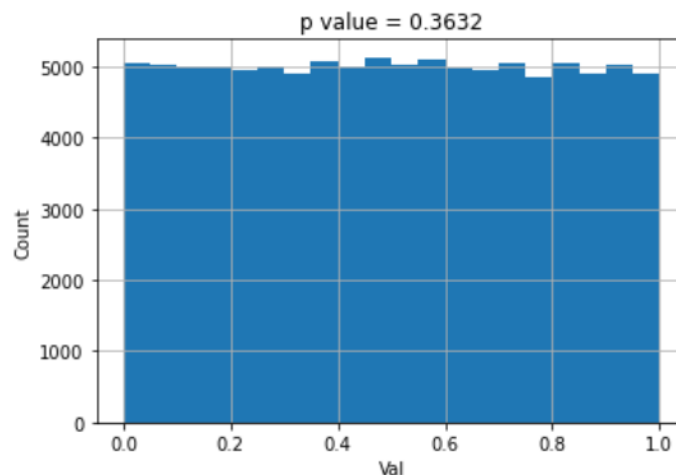
- Looking at the provided references a Chi-square test, to test for uniformity is very recommended, given it is mentioned in various of the references.

To do this I'll use python, so first I'll pass the .txt obtained from the C++ pseudo-random number generator, to a .csv, and then I'll use pandas and numpy to pass it into an array in Python, from which I can run the Chi-square test, as shown in class:

```

1
2 import scipy.stats as scp
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
6
7 df=pd.read_csv('C:\CARRERAS\MASTER UB\1r Semestre\Metodes Estadistics i
   matemantics\Exercices\Random Numbers\Random_numbers.csv', header =
   None, names=['random'])
8 x=df.to_numpy()
9
10 # Plot histogram; we use the hist function to recover the counts per bin
11 h= plt.hist(x,20) # h[0] will be the counts per bin
12 plt.grid(True)
13 plt.xlabel('Val')
14 plt.ylabel('Count')
15
16 # Run test
17 c2_stat, p_val = scp.chisquare(h[0])
18
19 t = plt.title("p value = {:.4f}".format(p_val))
20 plt.show()
21
22

```



from where we obtain p-values, for each of our 4 generators:

- Decent straight forward linear congruential: $p = 0.9367$
- Bad straight forward linear congruential (RANDU): $p = 0.2891$
- Minimal standard generator: $p = 0.3632$
- Rand() from C++: $p = 0.5869$

seeing that all the generators passed this test with a significance level of $\alpha = 0.25$ for each case, they also pass it for 5 more different seeds, which I tested.

- The next test we are going to do are tests of run above and below the median, also mentioned in various references as adequate ones, for this, we are going to add to our original C++ codes of each pseudo-random number generator this piece of code:

```

1
2 int counter=0;
3 for (int i=0; i<N ;i++){
4     if (ran[i]<0.5){
5         counter+=1;
6     }
7 }
8 printf("%d", counter);
9

```

For which we obtain the quantity of numbers that are in the wrong side of the expected median for each case, doing this for diverse initial seeds, and applying the statistic:

$$z = \frac{(x \pm 0.5) - \mu_x}{\sigma_x}$$

we need that $|z| < 2$ to pass the test with significance level $\alpha = 0.05$, and we obtain:

- Decent straight forward linear congruential: $z = -0.121$
- Bad straight forward linear congruential (RANDU): $z = 0.373$
- Minimal standard generator: $z = 0.157$
- Rand() from C++: $z = -0.105$

where we see that again, all test give satisfactory results, passing the test for significance level $\alpha = 0.05$, with great margin.

- Now, the following test is gonna be the Reverse Arrangements test, which will let us check if there is any correlation between the increase of the random generated numbers and the increase inside the computation time, or the opposite:

$$\sum_{i=1}^N \sum_{j=1(<i)}^N h_{ij} \quad \text{where } h_{ij} = 1 \text{ if } rand[i] > rand[j] \text{ or } 0 \text{ else}$$

this is implemented as:

```

1 int counter=0;
2 for (int i=0; i<N ;i++){
3     for (int j=0; j<N; j++){
4         if (i=j+1){
5             if (ran[i]<ran[j]){
6                 counter+=1;
7             }
8         }
9     }
10 }
11 printf("%d", counter);

```

for which we wanna obtain a result of $N/2$, for the generation to not be correlated (in this case $N=100000$). The obtained results are:

- Decent straight forward linear congruential: 49967

- Bad straight forward linear congruential (RANDU): 49983
- Minimal standard generator: 49918
- Rand() from C++: 49984

from where we see that all the number generators have good results, maybe the Minimal standard generator, have them a bit lower than the rest.

- And the final one will be the Overlapping Sums test, where we will sum the square difference deviation from the uniform distribution in each bin of the histogram:

$$\chi^2 = \sum_{i=1}^{Nbins} \frac{(O_i - E_i)^2}{E_i}$$

to compute this we will use:

```

1 counter=0
2 for i in range(len(h[0])):
3     print(h[0][i])
4     counter+=(h[0][i]-5000)**2/5000
5
6 print(counter)

```

For which we obtain, for each case a χ^2 :

- Decent straight forward linear congruential: $\chi^2 = 10.59$
- Bad straight forward linear congruential (RANDU): $\chi^2 = 21.90$
- Minimal standard generator: $\chi^2 = 20.53$
- Rand() from C++: $\chi^2 = 17.04$

checking the tables, for our $Nbins - 1 = 19$, and $\alpha = 0.05$ we get that χ must fulfill:

$$\chi^2 < 30.144$$

which all of our test fulfilled!

4 Discussion

The comparison, hasn't given much insight in which generators are better, maybe this is related with the fact, that I was using a laptop and the maximum length of the arrays of random numbers was of the order of 10^5 , which isn't near the periodicity of any of the generators method.

But even with the little differences, we can see in the Overlapping Sums test, that the decent straight forward linear congruential and the `rand()` from C++, where the two that performed better, when taking into account the square of the difference. Maybe if we did even more experiments, we could see a clear trend in favor of the "better" ones, but with the experiments we have used, it is not enough.

If we needed arrays of the order of 10^{10} random numbers, I imagine, the results would be pretty different, because in the bad ones, we would make the errors bigger each "cycle" when we repeat them, and the Overlapping Sums test, for example would notice that pretty quick.

But even if the test, wasn't strong enough, or we could have performed them better, the selection of this concrete test seem pretty good, we tested different things with each, that gave us an idea of how good the random generators are. The main ideas that we tested were:

- If the bins of a histogram follow a uniform distribution.
- If only if, half of the values are passed the 0.5 median value.
- If the values tend to get bigger or smaller in some period of the simulation.

References

- [1] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes in C*, Second. Cambridge, USA: Cambridge University Press, 1992.