

Enunciado 1 - Generación de Números y Variables Aleatorias

Joaquín Jiménez López de Castro — `jo.jimenez@alumnos.upm.es`

Alejandro Sánchez de Castro Fernández — `alejandro.sanchezdecastro.fernandez@alumnos.upm.es`

Ángel Fragua Baeza — `angel.fragua@alumnos.upm.es`

5 de octubre de 2021

1. Introducción

En esta práctica se pretende tomar contacto con algunas estrategias para probar la eficacia de generadores de números aleatorios. En concreto, se desea probar la eficacia de un generador congruencial multiplicativo de IMSL. Para ello, se utilizará la librería *TestU01* (Simard, 2013), que contiene herramientas para implementar y probar la eficacia de generadores de números aleatorios. Después se analizarán algunos de los contrastes empleados por la librería.

En la [Sección 2](#), se explicará como se ha implementado y probado el generador de IMSL en la librería *TestU01*. En la [Sección 3](#) se explicarán brevemente tres de los métodos que usa *TestU01* para comprobar la eficacia del generador. Finalmente en la [Sección 4](#) se harán unas conclusiones del trabajo realizado.

2. Implementación

Para la implementación de la práctica se ha usado una librería desarrollada en C llamada *TestU01*. En esta librería se encuentran implementados algunos de los generadores de números aleatorios más típicos junto con una serie de baterías de tests para comprobar si realmente son aleatorios o no. Para poder comenzar a usar esta librería se debe de seguir un proceso de descarga e instalación que se encuentra perfectamente explicado en su [guía](#).

Como para cualquier programa desarrollado en C, hemos generado un pequeño *Makefile* (véase Código 1) que nos permite agilizar el proceso de compilación y ejecución.

```
1 all: program
2
3 program:
4     gcc bat1.c -o bat1 -ltestu01 -lprobdist -lmylib -lm
5
6 run:
7 ifndef SEED
8     ./bat1
9 else
10    ./bat1 $(SEED) > results/seed-$(SEED).txt
11 endif
12
13 .EXPORT_ALL_VARIABLES:
14 LD_LIBRARY_PATH = /usr/local/lib
15 LIBRARY_PATH = /usr/local/lib
16 C_INCLUDE_PATH = /usr/local/include
```

Código 1: Makefile del programa

Dentro de este *Makefile* caben destacar ciertos detalles. En la línea 1.3 se establece la dependencia encargada de compilar el programa usando el típico programa *gcc* añadiéndole al final las librerías específicas de *TestU01*. En cuanto a la línea 1.6 incluye la dependencia para la ejecución del programa. Si esta dependencia se llama sin ningún argumento, ejecutará el programa haciendo uso de una semilla por defecto establecida en el programa, en cambio si se ejecuta acompañado de un argumento este se usará como semilla, guardando el output del programa en un archivo. Como la librería *TestU01* requiere de unas variables de sesión para la compilación, estas han sido añadidas al *Makefile* como un *target* tal y como se puede observar en la línea 1.13.

Todo el código necesario para llevar a cabo la generación de números aleatorios utilizando el generador congruencial multiplicativo de IMSL junto con el estudio de aleatoriedad se encuentra contenido en el archivo *bat1.c* tal y como se muestra en el Código 2.

```
1 #include <stdlib.h>
2 #include <stdio.h>
3 #include "ulcg.h"
4 #include "gdef.h"
5 #include "unif01.h"
6 #include "bbattery.h"
7
8 #define M 2147483647
9 #define A 16807
10 #define C 0
```

```
11 #define DEFAULT_SEED 559079219
12
13 int main (int argc, char* argv[])
14 {
15     int seed = (argc > 1) ? atoi(argv[1]) : DEFAULT_SEED;
16     unif01_Gen *gen;
17     gen = ulcg_CreateLCG (M, A, C, seed);
18     bbattery_SmallCrush (gen);
19     ulcg_DeleteGen (gen);
20
21     printf("Nuestra semilla: %d\n", seed);
22
23     return 0;
24 }
```

Código 2: bat1.c

2.1. Generador de números

En esta práctica se ha usado un generador congruencial, tratándose de uno de los generadores más populares. Este tipo de generadores fueron introducidos por Lehmer (1951), y siguen la función recursiva siguiente:

$$x_{n+1} = (ax_n + b) \mod m$$

Esta función genera números de forma aleatoria entre $[0, m)$, siendo a un multiplicador, "bün sesgo", m el módulo y x_0 la semilla inicial. En la práctica se solicitaba hacer uso del generador congruencial multiplicativo de IMSL siguiendo la siguiente fórmula:

$$x_{n+1} = (16807x_n) \mod (2^{31} - 1)$$

Para tener el código bien organizado se han declarado una serie de constantes. En la línea 2.8 se establece el valor del módulo m pero una vez hecha la operación $2^{31} - 1 = 2147483647$. En la línea 2.9 se establece el multiplicador a . En la línea 2.10 se establece el sesgo b , que como ya se ha mencionado ha de ser 0, ya que se trata de un generador congruencial multiplicativo. Finalmente, en la línea 2.11 se establece una semilla inicial (generada de forma aleatoria), para el caso de que el programa no reciba una semilla como argumento a la hora de ejecutarse.

Para crear el generador, la librería incluye una función llamada *ulcg_CreateLCG*, tal y como se puede observar en la línea 2.17, a la que se le pasan los parámetros típicos mencionados previamente. Esta función tiene la peculiaridad de que en vez de devolver los valores en el rango $[0, m)$ los devuelve en el rango $[0, 1)$ tal y como se nos solicitaba en el enunciado.

2.2. Batería de contrastes

Para este trabajo se va a ejecutar la batería de contrastes más sencilla, *bbattery_SmallCrush*. Cabe decir que algunos de los contrastes asumen que la muestra de números tiene al menos 30 bits de precisión, en caso de no ser así es altamente probable que fallen estos contrastes.

En la batería se aplican los siguientes contrastes:

1. smarsa_BirthdaySpacings con $N = 1, n = 5 \times 10^6, r = 0, d = 230, t = 2, p = 1$.
2. sknuth_Collision con $N = 1, n = 5 \times 10^6, r = 0, d = 216, t = 2$.
3. sknuth_Gap con $N = 1, n = 2 \times 10^5, r = 22, \alpha = 0, \beta = \frac{1}{256}$.
4. sknuth_SimpPoker con $N = 1, n = 4 \times 10^5, r = 24, d = 64, k = 64$.
5. sknuth_CouponCollector con $N = 1, n = 5 \times 10^5, r = 26, d = 16$.
6. sknuth_MaxOfn con $N = 1, n = 2 \times 10^6, r = 0, d = 105, t = 6$.
7. svaria_WeightDistrib con $N = 1, n = 2 \times 10^5, r = 27, k = 256, \alpha = 0, \beta = \frac{1}{8}$.
8. smarsa_MatrixRank con $N = 1, n = 20000, r = 20, s = 10, L = k = 60$.

9. `sstring_HammingIndep` con $N = 1, n = 5 \times 10^5, r = 20, s = 10, L = 300, d = 0$.

10. `swalk_RandomWalk1` con $N = 1, n = 10^6, r = 0, s = 30, L_0 = 150, L_1 = 150$.

Para comprobar la aleatoriedad del generador se han usado seis semillas para generar diferentes secuencias de números: 559079219, 4, 3, 2, 1 y 0. Se han sometido estas seis muestras de números a la batería de tests para comprobar el efecto de las semillas sobre el generador. Como era de esperar, la semilla 0 no pasa ningún contraste. En cuanto al resto de semillas hay tres contrastes que siempre fallan:

```

1 ===== Summary results of SmallCrush =====
2
3 Version:          TestU01 1.2.3
4 Generator:        ulcg_CreateLCG
5 Number of statistics: 15
6 Total CPU time:   00:00:05.37
7 The following tests gave p-values outside [0.001, 0.9990]:
8 (eps means a value < 1.0e-300):
9 (eps1 means a value < 1.0e-15):
10
11      Test                                p-value
12 -----
13 1 BirthdaySpacings                      eps
14 2 Collision                              eps
15 6 MaxOft                                eps
16 -----
17 All other tests were passed

```

Por lo que se entiende que el generador, por su naturaleza e independientemente de la semilla, no es capaz de pasar ninguno de estos contrastes.

3. Análisis de contrastes

3.1. Contraste Máximo-de-t

El contraste está descrito en Knuth (1988). Dada una muestra de números $U = \{U_1, U_2, \dots, U_n\}$, teniendo $U_i \in (0, 1)$ para cualquier i , sirve para rechazar la hipótesis de que U sigue una distribución uniforme en $(0, 1)$. El procedimiento consiste dividir U en $m = \lfloor \frac{n}{t} \rfloor$ clusters $C = \{C_1, \dots, C_m\}$, con $C_j = \{U_{jt}, U_{jt+1}, \dots, U_{j(t+1)-1}\}$ con $t \geq 1, j = 1, \dots, m$. Se obtiene $V = \{\max(C_1), \dots, \max(C_m)\} = \{V_1, \dots, V_m\}$. La hipótesis se rechaza si el test de Kolmogorov-Smirnov rechaza la hipótesis de que V tiene $F(x) = x^t, 0 \leq x \leq 1$ como función de distribución.

El motivo es que si U sigue una distribución uniforme, entonces:

$$\begin{aligned}
 F(x) &= P(V_j \leq x) = \\
 &= P(\max(\{U_{jt}, U_{jt+1}, \dots, U_{j(t+1)-1}\}) \leq x) = \\
 &= P(U_{jt} \leq x)P(U_{jt+1} \leq x) \dots P(U_{j(t+1)-1} \leq x) = x \dots x = x^t
 \end{aligned}$$

Este contraste admite variaciones. Una implícita es el parámetro t , donde si $t = 1$, se tiene una mera comprobación de que los valores de U son uniformes. En la librería *TestU01* se usa $t = 6$. También se puede intercambiar el test de Kolmogorov-Smirnov por el contraste χ^2 , que es lo que hace *TestU01*.

Es sencillo plantear una muestra que pase el test y tenga un patrón fácilmente observable, pero la utilidad de este contraste reside en rechazar generadores de números aleatorios, no en aceptarlos.

3.2. Collision

Este contraste también se encuentra en Knuth (1988). Para entender correctamente el test de *Collision*, podemos imaginarlo como un juego, en el que tenemos m urnas y n bolas que lanzamos de forma aleatoria. Para que este test tenga validez, se ha de cumplir que m es mucho mayor que n . Al ser m mucho mayor que n , lo más probable es que la mayoría de bolas entren en una urna vacía. Si se da el caso de que la urna ya tenía al menos una bola, se considera que se ha producido una colisión. El test de *Collision* cuenta

el número total de colisiones producidas, y si este no es ni demasiado grande ni demasiado pequeño el generador superará el test.

Partiendo de unos conocimientos básicos de estadística, podemos saber que la probabilidad de que una urna tenga k bolas es de:

$$p_k = \binom{n}{k} m^{-k} (1 - m^{-1})^{n-k}$$

Partiendo de la fórmula anterior, se puede calcular cuál es el número estimado de colisiones por urna.

$$\sum_{k \geq 1} (k-1)p_k = \sum_{k \geq 0} kp_k - \sum_{k \geq 1} p_k = \frac{n}{m} - 1 + p_0$$

Sabiendo cuál es el número estimado de colisiones por urna, se puede calcular fácilmente el número de colisiones totales estimadas. Para ello basta con multiplicar el número estimado de colisiones por urna, por el número total de urnas. Esta cuenta es algo costosa, puesto que la probabilidad de que una urna tenga 0 bolas $p_0 = (1 - m^{-1})^n = 1 - nm^{-1} + \binom{n}{2}m^{-2} - \text{términos menores}$ ya es de por sí costosa computacionalmente hablando. Por ello, para calcular el número total estimado de colisiones se suele usar la fórmula $\frac{n^2}{2m}$, la cual es menos costosa computacionalmente hablando y siempre da unos resultados ligeramente superiores a los exactos.

La potencia de este contraste se base en su capacidad para medir colisiones en múltiples dimensiones. Para ello basta con trocear la muestra inicial en vectores del mismo tamaño que número de dimensiones queramos usar en el test. Cada uno de estos vectores se usará para indexar una tabla de bits de tamaño m . Esta tabla se encontrará inicializada con todo ceros, y en el momento en el que se acceda a una posición cualquiera si hay un 0 se le transforma en 1, y si ya hay un 1 esto significa que ya se ha visitado esa posición, por lo que se mantiene el 1 y se incrementa en uno el número de colisiones.

En la librería *TestU01* la función encargada de llevar a cabo este contraste sería *sknuth_Collision*, cuya implementación se basa en otro test llamado *power divergence*. Este otro contraste está implementado por la función *smultin_Multinomial*, en donde uno de sus parámetros llamado *Sparse* permite que sus parámetros n y k sean muy grandes, lo cual hará que el número de urnas sea muy grande requisito principal del *Collision test*.

3.3. Birthday spacings

Este contraste fue propuesto por Marsaglia (1985). Se asume la premisa de que, en un conjunto de números ordenados escogidos de forma aleatoria, la distribución asintótica del número de valores duplicados entre los espacios de los números es Poisson con media $\lambda = \frac{m^3}{4n}$.

Es decir, se parte de una muestra $U = \{U_1, U_2, \dots, U_n\}$ de tamaño m , llamada cumpleaños. Los elementos de esta muestra tomarán valores entre $[0, n)$, siendo n el tamaño del espacio sobre el que se generan los elementos de la muestra; también se conoce a n como año. Esta muestra U se ordena de forma ascendente y se toman las diferencias entre los elementos consecutivos de forma $x_{j+1} - x_j$ tal que $1 < j < m$ de la muestra para calcular los espacios. Estos espacios se almacenan en otra muestra $K = \{K_1, K_2, \dots, K_n\}$ que vuelve a ordenarse para contar el número de colisiones, o lo que es lo mismo, la cantidad de valores duplicados. Esta cantidad de colisiones debe ajustarse al valor de λ . Este proceso se repite hasta obtener una muestra de 5000 valores de colisiones. Por último, se aplica una bondad de ajuste $\sum \frac{2(\text{observado} - \text{esperado})}{\text{esperado}}$ y ese valor se pasa a la función de distribución *chi-cuadrado* adecuada para obtener el p-valor y rechazar o no la hipótesis H_0 : “El número de colisiones es aproximadamente una variable aleatoria de Poisson con media igual a λ ”.

No hay estudios teóricos sólidos que encuentren m y n óptimos, pero sí hay generaciones de números aleatorios que cumplen esta regla de forma satisfactoria. Así como tampoco hay estudios concluyentes sobre el tamaño que debe tener n para poder comparar los resultados con la distribución de Poisson con media λ , pero la experiencia marca que para poder aplicar esta regla n debe ser al menos 2^{18} .

Para fallar este test el p-valor debe ser próximo a 0, pero fallar una vez es condición necesaria pero no suficiente. Para comprobar si un generador falla en este test es necesario probar con distintas semillas y ver si el p-valor difiere en órdenes de magnitud. Si un generador falla una vez con una semilla x_i con

un p-valor de $5x10^{-10}$ y con otra semilla x_j falla con $5x10^{-6}$ no puede afirmarse que el generador haya fallado el test. En cambio, si falla las dos veces con p-valor $5x10^{-10}$ puede afirmarse que el generador no ha pasado el test Mccullough (2006). Este es uno de los test más difíciles de pasar, si un conjunto pasa este test es muy probable que pase todos los demás.

En cuanto a la implementación, la librería *TestU01* ofrece la función *smarsa_BirthdaySpacings* para poder ejecutar esta prueba individualmente. La función incluye los siguientes parámetros:

$N \rightarrow$ Cantidad de veces que se repite el test sobre muestras de datos distintas del mismo conjunto de número aleatorios.

$n \rightarrow$ Número de “cumpleaños”. Tamaño de la muestra que se escoge del conjunto de números.

$d \rightarrow$ Tamaño del espacio.

$t \rightarrow$ Número de dimensiones del espacio.

$p \rightarrow$ Forma de ordenar el espacio. Puede tomar valores en $[1, 2]$.

$r \rightarrow$ Toma valores entre $[0, b]$ siendo b el número de bits que tiene el conjunto números aleatorios. Según este valor se ignoran los r bits más significantes del conjunto de número aleatorios.

En el caso de este trabajo se ha utilizado la batería de pruebas *bbattery_SmallCrush* en el que se ejecuta *smarsa_BirthdaySpacings* con unos valores concretos de sus parámetros. Estos valores son:

$N \rightarrow 1$

$n \rightarrow 5x10^{-10}$

$d \rightarrow 2^{30}$

$t \rightarrow 2$

$p \rightarrow 1$

$r \rightarrow 0$

4. Conclusión

En esta práctica se ha podido utilizar una librería de test de generadores de números aleatorios para comprobar la eficacia del generador congruencial multiplicativo de IMSL. Para ello, se ha pasado la batería de contrastes pequeña y tres tests han rechazado el generador como adecuado. Se han escogido estos tres contrastes para ser analizados y explicados brevemente.

Se concluye que el generador, aunque puede servir para algunos propósitos básicos, no es adecuado para funcionalidades más exigentes, como la criptografía.

Referencias

- Knuth, D. E. (1988). *Art of computer programming, volume 2: Seminumerical algorithms*. Addison-Wesley Professional, págs. 70-71.
- Lehmer, D. H. (1951). Mathematical Methods in Large-scale Computing Units. En: *Proceedings of the Second Symposium on Large Scale Digital Computing Machinery*. Cambridge, United Kingdom: Harvard University Press, págs. 141-146.
- Marsaglia, G. (1985). A current view of random number generators. En: *Computer Science and Statistics: The Interface*, págs. 3-10.
- Mccullough, B. D. (2006). A review of TESTU01. En: *Journal of Applied Econometrics* 21.5, págs. 677-682. DOI: [10.1002/jae.917](https://doi.org/10.1002/jae.917).
- Simard, R. J. (2013). A Software Library in ANSI C for Empirical Testing of Random Number Generators. En.