



DISTINTAS FORMAS DE SIMULAR VALORES DE VARIABLES ALEATORIAS CON DISTRIBUCIÓN NORMAL ESTÁNDAR

**DIFFERENT WAYS TO SIMULATE VALUES OF RANDOM
VARIABLES WITH STANDARD NORMAL DISTRIBUTION**

*Victor Miguel Ángel Burbano Pantoja**
*Margoth Adriana Valdivieso Miranda***
*Luis Alfonso Salcedo Plazas****

Recepción 20 /07/2011
Evaluación 10/08/2011
Aprobado 09/10/2011

Resumen

En este trabajo se presentan distintas formas de simular valores de variables aleatorias con distribución normal estándar utilizando los métodos propuestos por Papoulis (1991), Ríos (2000), Blanco (2004), Burbano (2010) y el simulador de SPSS. El objetivo es comparar los métodos

* Maestría, Docente Escuela de Matemáticas y Estadística UPTC, estudiante de Doctorado en Ciencias de la Educación RUDECOLOMBIA, CADE-TUNJA. UPTC. E-mail: victorburbanop@yahoo.es Trabajo realizado por el grupo Gamma de la UPTC, Tunja, proyecto “Simulación con modelos aleatorios y no aleatorios”.

** Maestría, Docente Escuela de Matemáticas y Estadística UPTC. E-mail: mavaldiviesom@yahoo.com

*** Especialista, Docente Escuela de Matemáticas y Estadística UPTC. E-mail: salcedop@email.com

Box-Muller, de las doce uniformes, un algoritmo presentado por Blanco en su libro “Probabilidad”, el simulador de SPSS y el método descrito por Burbano en su artículo publicado en la revista de la Facultad de Ciencias de la Universidad Pedagógica y Tecnológica de Colombia, el cual hace uso de la función inversa de la Distribución Lambda Generalizada citada por Karian y Dudewicz (2000).

Palabras clave: Simulación, modelos aleatorios, distribución normal estándar.

Abstract

In this paper we present different ways to simulate values of random variables with standard normal distribution using the methods proposed by Papoulis (1991), Ríos (2000), Blanco (2004), Burbano (2010) and SPSS simulator. The objective is to compare the methods: Box-Muller, the twelve uniforms, a Blanco’s algorithm showed in her book “Probabilidad”, the SPSS simulator and the method described by Burbano in his paper published in the Magazine of Sciences of Universidad Pedagógica y Tecnológica de Colombia, that paper use the inverse function of the Lambda Generalized Distribution named by Karian y Dudewicz (2000).

Keywords: Simulation, random models, standard normal distribution.



Introducción

La simulación por computador tiene actualmente una enorme aplicación en diversos campos del conocimiento humano, tales como: estadística, economía, física, informática, ingeniería, entre otros. La simulación se está utilizando para resolver gran variedad de problemas que por métodos analíticos no se habían podido solucionar. Autores dedicados a este campo indican que para hacer procesos de simulación es conveniente partir de modelos matemáticos que permitan de manera razonable emular el comportamiento de un fenómeno o de un sistema de interés. Para implementar dichos modelos se requiere usar elementos teóricos provenientes de la estadística matemática, la teoría de probabilidad, las ecuaciones diferenciales y la programación de computadores, entre otras.

En el campo de la estadística, la distribución normal ha sido considerada como una pieza fundamental en muchos procesos de inferencia, en los cuales es necesario que se cumpla el supuesto de normalidad. La distribución normal también ha permitido modelar una gran cantidad de fenómenos que ocurren con cierta regularidad en la naturaleza. Por lo anterior, la distribución normal (estándar) ha sido estudiada ampliamente en forma analítica y mediante simulaciones, constituyéndose en muchos casos en el límite cuando se trabaja teoría asintótica de variables aleatorias.

Métodos

A continuación se describen brevemente cinco diversos métodos utilizados para simular valores de variables aleatorias con distribución normal estándar: el método de Box-Muller mencionado en Papoulis (1991), el de las doce uniformes descrito por Ríos (2000), el presentado por Blanco (2004) basado en el método del rechazo, el de la transformada inversa usando la Distribución Lambda Generalizada

propuesto por Burbano (2010) y el del generador de valores aleatorios normales estándar de SPSS.

Método de Box-Muller

Según Papoulis (1991), para aplicar el método de Box-Muller es conveniente considerar dos variables aleatorias independientes U_1, U_2 , cada una con distribución uniforme en el intervalo $(0,1)$ que permiten generar los valores u_1, u_2 correspondientes a las variables mencionadas anteriormente. Para generar los valores z_1, z_2 correspondientes a las variables aleatorias independientes Z_1 y Z_2 , cada una con distribución normal estándar, es posible utilizar las expresiones siguientes:

$$z_1 = \sqrt{-2\text{Ln}(u_1)} \cos(2\pi u_2) \quad (1)$$

$$z_2 = \sqrt{-2\text{Ln}(u_1)} \text{sen}(2\pi u_2) \quad (2)$$

Método de las doce uniformes

Para Ríos (2000), se consideran n variables aleatorias independientes U_1, U_2, \dots, U_n , cada una con distribución uniforme en el intervalo $(0,1)$ que permiten generar los valores u_1, u_2, \dots, u_n correspondientes a las variables ya mencionadas. Para generar un valor z correspondiente a la variable aleatoria Z con distribución normal estándar, se utiliza la siguiente expresión:

$$z = \frac{\sum_{i=1}^n u_i - \frac{n}{2}}{\sqrt{\frac{n}{12}}}$$

Entre más grande sea n , concordante con el teorema central del límite, mayor certeza se tendrá de que z corresponda a un valor proveniente de una distribución normal estándar. Es suficiente tomar $n=12$, para que el valor z resultante



se considere que pertenece a una distribución normal estándar, admitiendo un error que se estima pequeño. Así, reemplazando en la expresión anterior, se tiene:

$$z = \frac{\sum_{i=1}^{12} u_i - \frac{12}{2}}{\sqrt{\frac{12}{12}}} = \sum_{i=1}^{12} u_i - 6$$

Luego,

$$z = \sum_{i=1}^{12} u_i - 6 \quad (3)$$

La expresión (3) se conoce como el método de las doce uniformes.

Método del rechazo para simular valores de una distribución normal

De acuerdo con Blanco (2004), en general, el método del rechazo supone que para generar valores x de una variable aleatoria X con función de densidad de probabilidad f , se puede generar una variable aleatoria Y con una función de densidad g , luego se genera un número aleatorio RND perteneciente al intervalo $(0,1)$ y se acepta este valor generado con una probabilidad proporcional a:

$$\frac{f(Y)}{g(Y)}$$

Luego, se aplica el siguiente algoritmo:

Se genera Y con densidad g .
Se genera un número aleatorio RND .

Si $RND \leq \frac{f(Y)}{cg(Y)}$, entonces hacer $X = Y$ en caso contrario regresar a 1.

La variable aleatoria X tiene función de densidad f .

Para el caso particular de que X corresponda a una variable aleatoria con distribución normal estándar, se tiene que su función de densidad de probabilidad es:

$$f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$$

$x \in \mathbb{R}$, para generar los valores de X se observa en primer lugar que la variable aleatoria $W = |X|$ tiene como función de densidad

$$h(x) = \frac{2}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$$

$x \in \mathbb{R}$, Posteriormente, se aplica el algoritmo siguiente: se genera primero una variable aleatoria W , se obtienen luego los valores de la variable X al hacer que sea igualmente probable que W sea igual a X ó a $-X$.

Para generar W se usa el método del rechazo con $g(x) = \exp(-x)$ para $x > 0$, se observa que la función $\frac{h(x)}{g(x)}$ toma su valor máximo en $x = 1$.

Luego,

$$c = \frac{2}{\sqrt{2\pi}} \exp\left(\frac{1}{2}\right)$$

Por lo tanto,

$$\frac{h(x)}{cg(x)} = \exp\left(-\frac{1}{2}(x-1)^2\right)$$

El algoritmo se resume en los siguientes pasos: (4)

Entrada: número máx de simulaciones.

Salida: valores z correspondientes a la variable aleatoria Z .



Para i , desde 1 hasta $máx$, haga:

Inicialización:

Genere una variable aleatoria Y con distribución exponencial de parámetro $\lambda = 1$

Genere un número aleatorio RND

Verificación:

Mientras $RND > \exp\left(-\frac{1}{2}(Y-1)^2\right)$ haga:

Genere una variable aleatoria Y con distribución exponencial de parámetro $\lambda = 1$

Genere un número aleatorio RND

Salida: (para escribir Y ó $-Y$ con igual probabilidad)

Genere un número aleatorio rnd

Si $rnd > 0.5$ entonces escriba $-Y$

En caso contrario escriba Y

Método de la transformada inversa usando la Distribución Lambda Generalizada (DLG)

Burbano (2010), en su artículo publicado en la revista de la Facultad de Ciencias de la Universidad Pedagógica y Tecnológica de Colombia, propone utilizar la función inversa de la Distribución Lambda Generalizada citada por Karian y Dudewicz (2000), para generar valores de variables aleatorias con distribución normal estándar con base en el método de la transformada inversa. El método se resume en lo siguiente:

En la Distribución Lambda Generalizada definida por la siguiente expresión:

$$F^{-1}(y) = F^{-1}(y, \lambda_1, \lambda_2, \lambda_3, \lambda_4) = \lambda_1 + \frac{y^{\lambda_3} - (1 - y^{\lambda_4})}{\lambda_2}$$

con $\lambda_2 \neq 0$, $0 \leq y \leq 1$ donde λ_1 es el parámetro de localización, λ_2 es el parámetro de escala, λ_3 determina el

sesgo (coeficiente de asimetría) y λ_4 determina la curtosis. La distribución normal estándar se obtiene asignando los siguientes valores específicos a sus parámetros:

$$\lambda_1 = 0, \lambda_2 = 0.1975, \lambda_3 = 0.1349, \lambda_4 = 0.1349$$

resultando

$$F^{-1}(y) = \frac{y^{0.1349} - (1-y)^{0.1349}}{0.1975} = x$$

Los valores z de la variable aleatoria Z con distribución normal estándar se obtienen generando valores u correspondiente a una variable aleatoria U con distribución uniforme en el intervalo $(0,1)$ que remplazan a los valores “ y ” en la anterior expresión, obteniéndose:

$$z = F^{-1}(u) = \frac{u^{0.1349} - (1-u)^{0.1349}}{0.1975} \quad (5)$$

2.5 Valores de una variable aleatoria normal estándar usando el generador de SPSS

Otra forma de generar valores z de la variable aleatoria Z con distribución normal estándar es utilizando la siguiente función disponible en el paquete estadístico SPSS,

$$z = \text{RV. NORMAL}(0,1) \quad (6)$$

Para acceder a la anterior función, se ubica la opción *Transformar* del menú principal de SPSS, de ella se selecciona la opción *Computar*, y del grupo de funciones se escoge la función *Random numbers*, y de esta se elige la que corresponde a *Rv.Normal*, escribiendo entre paréntesis $(0,1)$, lo cual indica que se trabajará con una media igual a cero y una desviación estándar igual a 1 sobre la distribución normal.



Resultados

Obtención de valores simulados

A continuación se presentan los resultados obtenidos al generar 20 valores z de la variable aleatoria Z con distribución normal estándar, mediante los cinco métodos indicados cuyas expresiones corresponden a (1), (2), (3), (4), (5) y (6), a fin de comparar cuál de ellos se ajusta mejor a una distribución normal estándar. Para determinar si los valores simulados para cada caso se adaptaban a una distribución normal, se utilizaron dos criterios: la prueba K-S de Kolmogorov- Simirnov usando un nivel de significancia del 5%, y el gráfico P-P, que corresponde a las probabilidades empíricas comparadas con las probabilidades teóricas de dichas observaciones sobre una distribución normal estándar. Entre más cerca estén los puntos, mejor será su ajuste a la distribución teórica ya mencionada.

Tabla 1. Valores simulados con cada uno de los cinco métodos

Box-Muller	Doce-Uniform	Blanco-Rech	DLG-Burbano	SPSS-Norm
1,503366	-0,742584	-0,127618	-1,17391	1,36357
2,218829	-0,629593	0,0607	-1,56646	-0,21403
0,694112	1,623618	-0,208279	-0,881692	-0,31856
1,741293	-0,837836	1,126611	-1,28303	0,69544
-0,397529	-1,14753	-1,17977	0,45326	-0,9808
-0,302155	0,348042	2,015836	0,500258	-0,00344
0,358106	-0,760116	0,432679	1,108846	1,0578
-0,859316	0,295364	0,686019	-0,379449	-1,38483
-1,18319	-0,688431	1,210588	-0,0089	-0,28052
-0,249991	0,514155	-1,48453	0,526872	-0,77964
1,408396	0,816556	0,069264	-0,816715	-0,31455
0,860712	-0,004944	-0,748818	0,746395	-1,13914
1,691339	-1,6346	1,084233	-1,50059	1,14743
1,25775	-1,33495	-0,379342	0,067423	1,35468

-0,790872	-0,088113	0,73966	0,414752	-1,22224
-0,801971	-1,38841	1,451793	1,522194	-0,03365
-0,397083	-0,658262	0,135903	-0,476765	0,39979
1,1637	0,279496	-1,10142	0,056575	0,0264
0,026586	0,659167	-0,412945	0,721593	0,40659
-0,80325	-1,10495	-2,32064	-1,13804	0,12717

A fin de tener un contexto que permita hacer una comparación con el menor sesgo posible, se generaron 20 números aleatorios con distribución uniforme en el intervalo (0,1), y con base en ellos se obtuvieron los valores que aparecen en la tabla 1, para cada caso en una sola corrida, de la siguiente forma: se utilizó la expresión (5) directamente, para el método de Box-Muller se utilizaron los diez primeros números aleatorios en u_1 y los diez siguientes en u_2 en las expresiones (1) y (2); para el método de las doce uniformes se seleccionaron 20 muestras aleatorias de tamaño doce del conjunto de los 20 números aleatorios, y para el método del rechazo fue necesario generar algunos números aleatorios adicionales a fin de completar el algoritmo y obtener los 20 valores que se indican en la anterior tabla. En estas circunstancias es posible que se haya generado un contexto de condiciones favorables para el método del rechazo y desfavorable para el método de las doce uniformes.

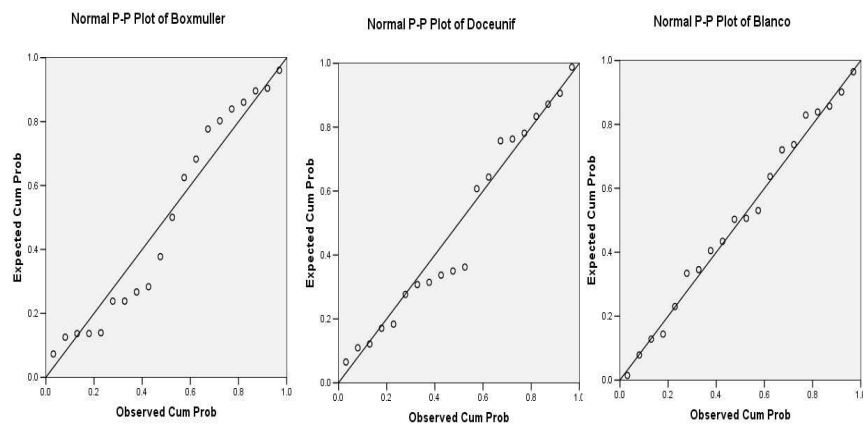
En la tabla 2 se presentan los resultados del análisis exploratorio de datos sobre los cinco grupos de datos correspondientes a los valores simulados con cada método, obteniéndose la media y la desviación estándar, además se presenta para cada caso el P-valor obtenido de la aplicación de la prueba de K-S de Kolmogorov-Smirnov y que servirá de criterio para determinar cuál método permite obtener los datos simulados que mejor se ajusten a una distribución normal estándar.



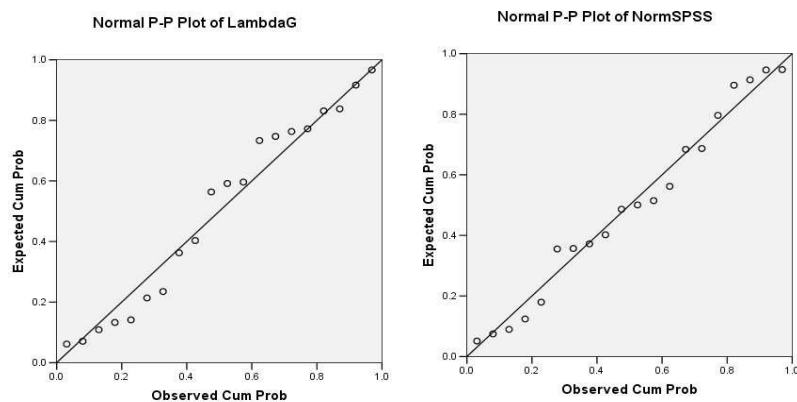
**DISTINTAS FORMAS DE SIMULAR VALORES DE
VARIABLES ALEATORIAS CON DISTRIBUCIÓN
NORMAL ESTÁNDAR**

	Box-Muller	12-Uniform	Recha Blanco	DLG-Burbano	SPSS
Media	0,35694137	-0,3241956	0,0524962	-0,15536919	-0,004626
Desviac	1,0593167	0,86597064	1,085488	0,914658247	0,8448104
P-valor	0,635	0,481	0,999	0,868	0,980

Tabla 2. Valores de la media, desviación estándar y P-valor para la prueba K-S



Gráfica 1. Gráficos P-P para el método de Box-Muller, doce uniformes y del rechazo (Blanco, 2004)



Gráfica 2. Gráficos P-P para el método Lambda Generalizada y el generador de SPSS

En las gráficas 1 y 2 se presentan los gráficos P-P correspondientes a los métodos: Box-Muller, doce uniformes, del rechazo (Blanco), de la transformada inversa con la familia Lambda Generalizada y el generador de SPSS.

Discusión de resultados

De la tabla 2, concordante con un análisis exploratorio de datos, se observa e interpreta que el método que proporcionó un valor para la media más cercano a cero fue el generador de SPSS (-0,0046269), le siguen el método del rechazo y el basado en la Distribución Lambda Generalizada, y el que más se alejó de cero fue el método de Box-Muller (0,35694137). El método que proporcionó un valor de la desviación estándar más cercano a 1 fue el de Box-Muller (1,0593167), seguido del método del rechazo y el basado en la Distribución Lambda Generalizada, y el que más se alejó de 1 fue el generador de SPSS (0,84481042).

Las gráficas 1 y 2 permiten visualizar que los puntos correspondientes al método del rechazo son los que se encuentran más cerca de la recta diagonal, proporcionando una idea gráfica de que este método es el que mejor ajusta los datos simulados a una distribución normal estándar. Le siguen el generador de SPSS y el basado en la Distribución Lambda Generalizada, y los puntos más distantes corresponden al método de las doce uniformes. Lo anterior permite afirmar que todos los métodos sí sirven para simular valores de una variable aleatoria con distribución normal, aunque dejan una primera sensación de que el método del rechazo es el más eficiente y el de las doce uniformes el menos eficiente. El método basado en la Distribución Lambda Generalizada se desempeña bastante bien, pero además tiene la ventaja de ser muy versátil puesto que permitiría generar valores aleatorios hasta con una calculadora corriente.



Los P- valor que se indican en la tabla 2 corresponden a pruebas de hipótesis como las que se indican a continuación, para cada uno de los cinco ajustes considerando un nivel de significancia del 0.05:

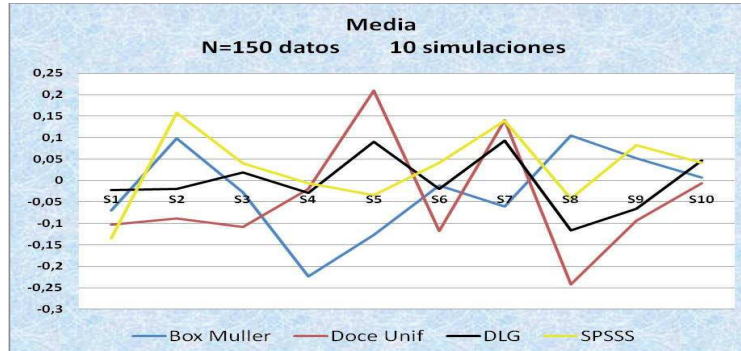
H_0 : los valores simulados provienen de una distribución normal estándar.

H_1 : los valores simulados no provienen de una distribución normal estándar.

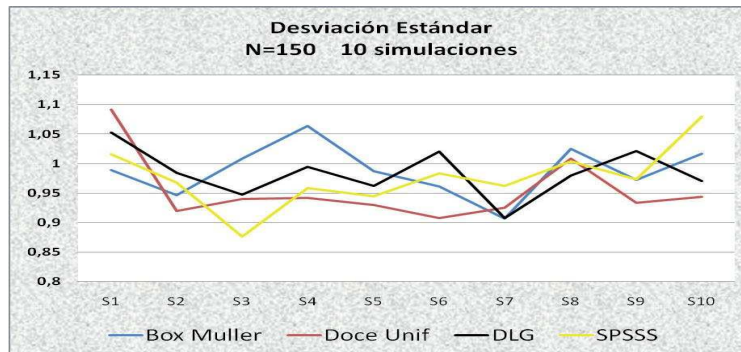
Los P-valores mencionados permiten concluir que en los cinco métodos los valores simulados sí provienen de una distribución normal estándar debido a que el P-valor en cada caso no es menor que 0.05. Sin embargo, el P-valor más grande lo presenta el método del rechazo (Blanco, 2004), le siguen el generador de SPSS y el basado en la Distribución Lambda, y el que menor P-valor presenta es el método de las doce uniformes. Es posible que la disminución de su eficiencia se deba al muestreo que fue necesario realizar del conjunto de 20 números aleatorios tomados como base para efectuar la comparación de los citados métodos.

Mayor número de simulaciones

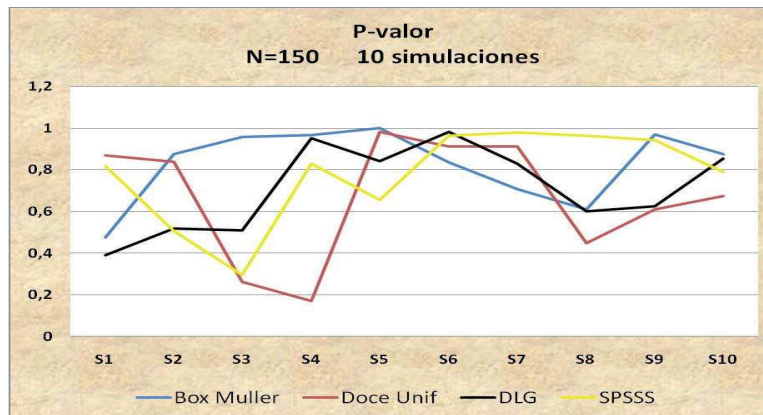
Nuevamente se generaron 150 valores z de la variable aleatoria Z con distribución normal estándar, mediante los cinco métodos indicados cuyas expresiones corresponden a (1), (2), (3), (4), (5) y (6), tomando como base para los cinco métodos, 150 números aleatorios con distribución uniforme en el intervalo (0,1). El procedimiento se repite diez veces (se hacen diez simulaciones) cada caso y se obtienen los valores de la media, la desviación estándar y el P-valor (no se presentan aquí las tablas de datos por razones obvias), pero con ellos se han elaborado las gráficas correspondientes al comportamiento de la media para los distintos métodos (gráfica 3), al comportamiento de la desviación estándar (gráfica 4) y para los P-valor (gráfica 5).



Gráfica 3. Comportamiento de la media para las diez simulaciones S_i , $i=1,2,\dots,10$



Gráfica 4. Comportamiento de la desviación estándar para las diez simulaciones S_i , $i=1,2,\dots,10$



Gráfica 5. Comportamiento del P-valor para las diez simulaciones S_i , $i=1,2,\dots,10$



Al observar en la gráfica 3 las líneas quebradas, se deduce que el método correspondiente a la Distribución Lambda Generalizada (línea quebrada negra) es el que más cerca queda de la línea horizontal correspondiente a la media igual a cero; le siguen en su orden el generador de SPSS (la línea amarilla), el método de las doce uniformes y finalmente el método de Box-Muller, no se ha incluido el método del rechazo porque no genera comparabilidad por la estructura del método al tomar como base los 150 números aleatorios.

En la gráfica 4, observando las líneas quebradas, se deduce que todos los métodos presentan una desviación estándar cercana a 1 (toman valores en el intervalo (0.95, 1.05)) mostrando poca variabilidad. Sin embargo, el método basado en la Distribución Lambda Generalizada está más próximo a la línea horizontal correspondiente a la desviación estándar igual a uno, le siguen en su orden el generador de SPSS, el método de Box-Muller y finalmente el de las doce uniformes, que presenta leve mayor variabilidad.

Observando las líneas quebradas de la gráfica 5, se deduce que todos los métodos presentan un P-valor mayor que 0.05, lo cual indica que todos los métodos ajustan los valores simulados a una distribución normal estándar. Sin embargo, hay un predominio de valores cercanos a 1 en el método de Box-Muller hasta la simulación cinco, seguida del método basado en la Distribución Lambda Generalizada, el generador de SPSS y finalmente el método de las doce uniformes. De la simulación seis a la simulación diez, mejora el simulador de SPSS, seguido del método basado en la Distribución Lambda Generalizada, el método de las doce uniformes y, finalmente, el método de Box-Muller disminuye su eficiencia.

Conclusiones

Existen diversos métodos para generar valores de una variable aleatoria con distribución normal estándar. Aquí se han originado valores de una variable aleatoria con distribución normal estándar por cinco métodos que permiten concluir que todos ajustan a la mencionada distribución, porque su P-valor es mayor que 0.05. Sin embargo, unos resultan más eficientes que otros, tomando el P-valor como medida de eficiencia.

Para los 20 números aleatorios generados, el método del rechazo resultó con mayor P-valor de 0.999, le siguen el simulador de SPSS, el método basado en la DLG, el método de Box-Muller y, finalmente, el de las doce uniformes.

Para 150 números aleatorios generados y en diez simulaciones, se encuentra que el método de las doce uniformes es el que tiene mayor desvío en cuanto a la media y al P-valor. Además, se observa que el método de SPSS es el que tiene mayor alejamiento en cuanto a la desviación estándar.

El método de la transformada inversa posibilita la utilización directa de la Distribución Lambda Generalizada con su función percentil, para generar de manera eficiente valores de variables aleatorias con distribución normal estándar y se constituye en una forma alterna a los métodos clásicos de simular variables aleatorias continuas.

Lista de referencias

- Azarang, M. R. (1996). *Simulación y análisis de modelos estocásticos*. Mexico: McGrawHill.
- Blanco, L. (2004). *Probabilidad*. Bogotá: Universidad Nacional de Colombia.



DISTINTAS FORMAS DE SIMULAR VALORES DE
VARIABLES ALEATORIAS CON DISTRIBUCIÓN
NORMAL ESTÁNDAR

- Karian, Z.A. & Dudewicz, E.J. (2000). *Fitting Statistical Distributions: The Generalized Lambda Distribution and Generalized Bootstrap Methods*. Boca Ratón, FL.: CRC Press.
- Leiva, V., Sanhueza, A., Sen, P. & Paula, G. (2008). *Journal of Statistical Computation and Simulation* 78, (11), 1105-1118.
- Papoulis, A. (1991). *Probability, Random Variables and Stochastic Process*. New York: McGraw-Hill.
- Ríos, D., Ríos, S. & Jiménez, J. M. (2000). *Simulación, métodos y aplicaciones*. Bogotá: Alfaomega.
- Ross, S. (1998). *A First Course in Probability*. United States of America: Prentice Hall.

