

EST-24107: Simulación

Profesor: Alfredo Garbuno Iñigo — Primavera, 2022 — Números no uniformes.

Objetivo: En esta sección del curso veremos métodos para generar números aleatorios de distribuciones un poco mas generales que una distribución uniforme. En particular veremos algunos ejemplos de distribuciones continuas, discretas y mezclas. Además veremos una prueba de uniformidad adicional a la prueba Kolmogorov-Smirnov.

Lectura recomendada: Capítulo 2 de [2].

1. GENERACIÓN DE VARIABLES NO UNIFORMES

R, por ejemplo, tiene distintos generadores de variables aleatorias. Un catálogo nos muestra que podemos generar mucho mas variables que solamente una variable uniforme.

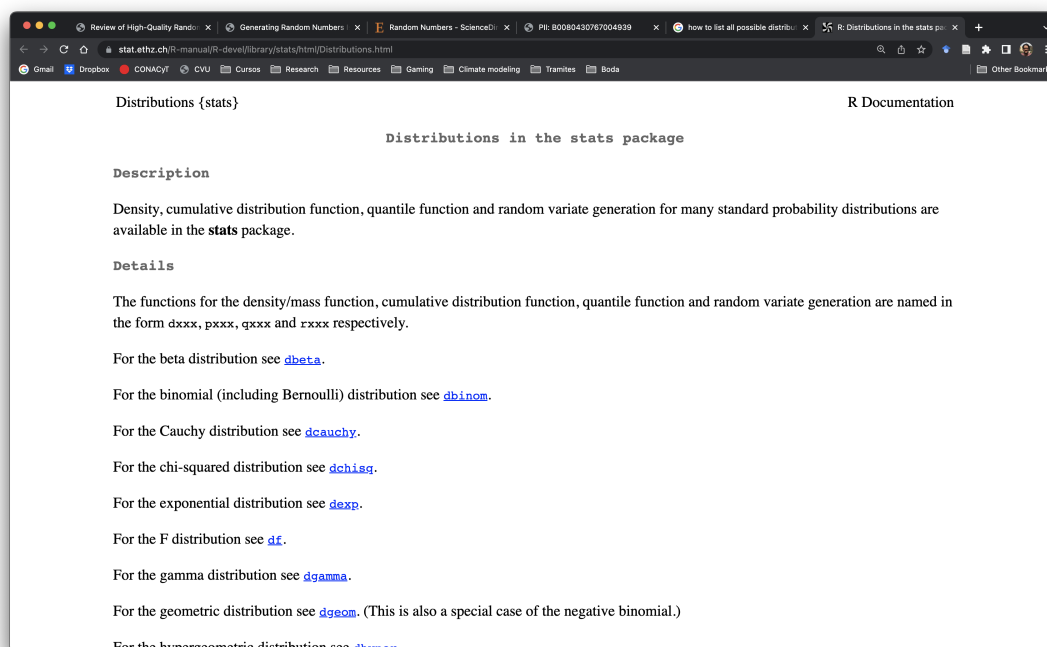


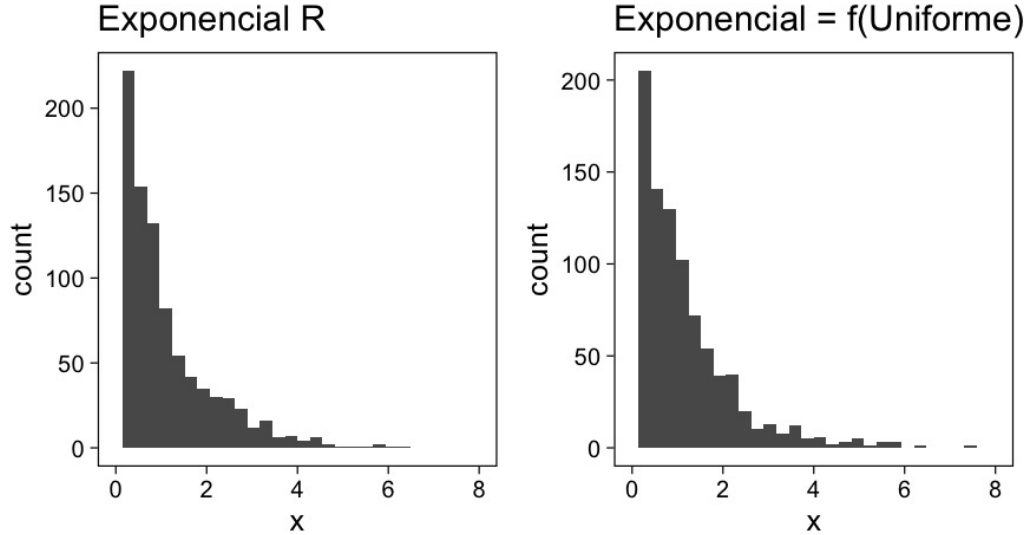
FIGURA 1. Catálogo de distribuciones en R.

En esta sección exploraremos los mecanismos tradicionales para generar números pseudo-aleatorios de distribuciones un poco mas generales que una distribución uniforme.

1.1. Método de la transformada inversa

De su curso de cálculo de probabilidades se acordarán que si tenemos una variable aleatoria X y una función de acumulación \mathbb{P} . Si utilizamos $U = \mathbb{P}(X)$, con U una variable aleatoria con distribución $U(0, 1)$. Entonces, $X \sim \mathbb{P}$.

1.1.1. Ejercicio: Considera $X \sim \text{Exp}(1)$, de tal forma que $\mathbb{P}(x) = 1 - e^{-x}$. Si sólo sabemos generar números uniformes, ¿cómo generarías números de X ?



1.1.2. *Tarea:* Considera las siguientes distribuciones:

1. Logística con μ, β y función de acumulación

$$\mathbb{P}(x) = \frac{1}{1 + e^{-(x-\mu)/\beta}}. \quad (1)$$

2. Cauchy con parámetros μ, β y función de acumulación

$$\mathbb{P}(x) = \frac{1}{2} + \frac{1}{\pi} \arctan((x - \mu)/\beta). \quad (2)$$

Genera números aleatorios con valores $\mu = 1$, $\beta = 2$ y comenta el rol de cada parámetro en cada distribución.

1.2. Generalización

En el caso anterior asumimos que la función inversa de la función de acumulación existe. Sin embargo, no siempre es el caso. Necesitamos definir lo siguiente.

1.2.1. *Definición [Función inversa generalizada]:* Sea X una variable aleatoria con función de acumulación \mathbb{P} . La función inversa generalizada está definida como

$$\mathbb{P}^{-1}(u) = \inf\{x | F(x) \geq u\}. \quad (3)$$

1.2.2. *Teorema [de la función inversa generalizada]:* Sea \mathbb{P} una función de distribución. Sea $\mathbb{P}^{-1}(\cdot)$ la función inversa generalizada de \mathbb{P} y $U \sim U(0, 1)$. Entonces, la variable aleatoria $X = \mathbb{P}^{-1}(U)$ tiene distribución \mathbb{P} .

2. DISTRIBUCIONES CONTINUAS

Hay situaciones donde generar números aleatorios puede ser extremadamente sencillo pues existe una relación entre la distribución que queremos generar y ciertos componentes fáciles de simular.

2.1. Variables χ^2

Recordarán de su curso de cálculo de probabilidades que si $X_i \sim \text{Exp}(1)$ entonces

$$Y = 2 \sum_{j=1}^{\nu} X_j \sim \chi_{2\nu}^2, \quad (4)$$

con $\nu \in \mathbb{N}$.

2.2. Variables gamma

También recordarán que si $X_i \sim \text{Exp}(1)$ entonces

$$Y = \beta \sum_{j=1}^{\alpha} X_j \sim \text{Gamma}(\alpha, \beta), \quad (5)$$

con $\alpha \in \mathbb{N}$.

2.3. Variables beta

También recordarán que si $X_i \sim \text{Exp}(1)$ entonces

$$Y = \frac{\sum_{j=1}^{\alpha} X_j}{\sum_{j=1}^{\alpha+\beta} X_j} \sim \text{Beta}(\alpha, \beta), \quad (6)$$

con $\alpha, \beta \in \mathbb{N}$.

2.3.1. Tarea: Prueba las igualdades anteriores.

2.3.2. Código: Por ejemplo, podemos utilizar el siguiente código para generar variables de una χ_6^2 .

```
1 set.seed(108)
2 U <- runif(3 * 10^4)      # Genera uniformes
3 U <- matrix(U, nrow = 3)  # Transforma a matriz
4 X <- -log(U)              # Transforma a exponenciales
5 X <- 2 * apply(X, 2, sum)  # Suma los tres renglones
6 summary(X)
```

```
1  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
2 0.151  3.505   5.420   6.064  7.926  27.904
```

A partir de la versión 4.1.1 R cuenta con un operador especial (`|>`) llamado **pipe** el cual permite *anidar* ciertas funciones y evitar la asignación repetitiva de variables.

```
1 set.seed(108)
2 runif(3 * 10^4) >        # Genera uniformes
3   matrix(nrow = 3) >     # Transforma a matriz
4   log() >                # Calcula logaritmos
5   apply(2, function(x){-2 * sum(x)} ) >
6   summary()
```

```
1  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
2 0.151  3.505   5.420   6.064  7.926  27.904
```

2.4. Variables Gaussianas correlacionadas

Supongamos que queremos generar un par de variables $X \in \mathbb{R}^2$ de tal forma que

$$X \sim N(0, \Sigma), \quad (7)$$

donde $\Sigma_{ii} = 1$ para $i \in \{1, 2\}$ y $\Sigma_{ij} = \rho$ con $i \neq j$. Supongamos que sólo sabemos generar números aleatorios $N(0, 1)$.

¿Cómo generaríamos los vectores aleatorios que necesitamos?

¿Qué saben de propiedades matriciales de vectores aleatorios?

```

1 set.seed(108)
2 Sigma <- diag(2); Sigma[1,2] <- .75; Sigma[2,1] <- .75;
3 L <- chol(Sigma)
4
5 Z <- rnorm(2 * 10^4)      # Generamos vectores estandar
6 Z <- matrix(Z, nrow = 2)  # Reacomodamos en matriz
7 X <- t(L) %*% Z           # Transformacion lineal
8 cov(t(X))

```

```

1      [,1]      [,2]
2 [1,] 1.0173 0.7772
3 [2,] 0.7772 1.0312

```

El operador `%*%` ejemplifica uno de los limitantes por diseño de **R**. Pues no está hecho para realizar operaciones vectoriales de manera nativa. Por ejemplo, en **Matlab** las operaciones son nativas y en **python** a través de **numpy** las operaciones matriciales también (y parte de los métodos).

3. DISTRIBUCIONES DISCRETAS

Ahora, veremos algunas técnicas generales para distribuciones discretas. O mejor dicho, para generar números aleatorios con soporte en los enteros.

Supongamos que nuestro objetivo es poder generar de una $X \sim \mathbb{P}_\theta$ donde $X \in \mathbb{Z}$. La estrategia es guardar todas las probabilidades del soporte. Es decir, calcular

$$p_0 = \mathbb{P}_\theta(X \leq 0), \quad p_1 = \mathbb{P}_\theta(X \leq 1), \dots, \quad (8)$$

generar $U \sim U(0, 1)$ y establecer

$$X = k \text{ si } p_{k-1} < U < p_k. \quad (9)$$

3.1. Binomial

Supongamos que nos interesa $X \sim \text{Bin}(10, 0.3)$, el vector de probabilidades lo podemos calcular con la función `pbinom(k, 10, .3)`.

```

1 k <- 1:10
2 pbinom(k, 10, .3)

```

```

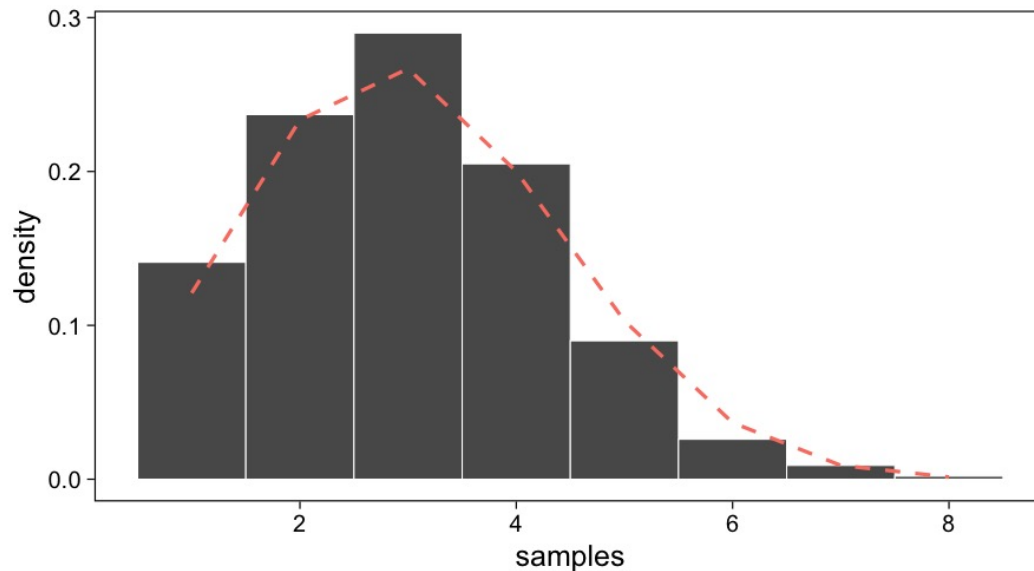
1 [1] 0.1493 0.3828 0.6496 0.8497 0.9527 0.9894 0.9984 0.9999 1.0000 1.0000

```

```

1 rbinomial <- function(nsamples, size, theta){
2   probs <- pbinom(1:10, size, theta)
3   x <- c()
4   for (jj in 1:nsamples){
5     u <- runif(1)
6     x[jj] <- which(probs > u)[1]
7   }
8   return(x)
9 }

```



3.2. Poisson

Ahora supongamos que nos interesa simular de una Poisson con parámetro $\lambda = 7$.

¿Cuál es el soporte de una $\text{Bin}(10, .3)$? ¿Cuál es el soporte de una $\text{Poisson}(7)$?

Tenemos que guardar las probabilidades

```

1 k <- 1:24
2 ppois(k, 7)

```

```

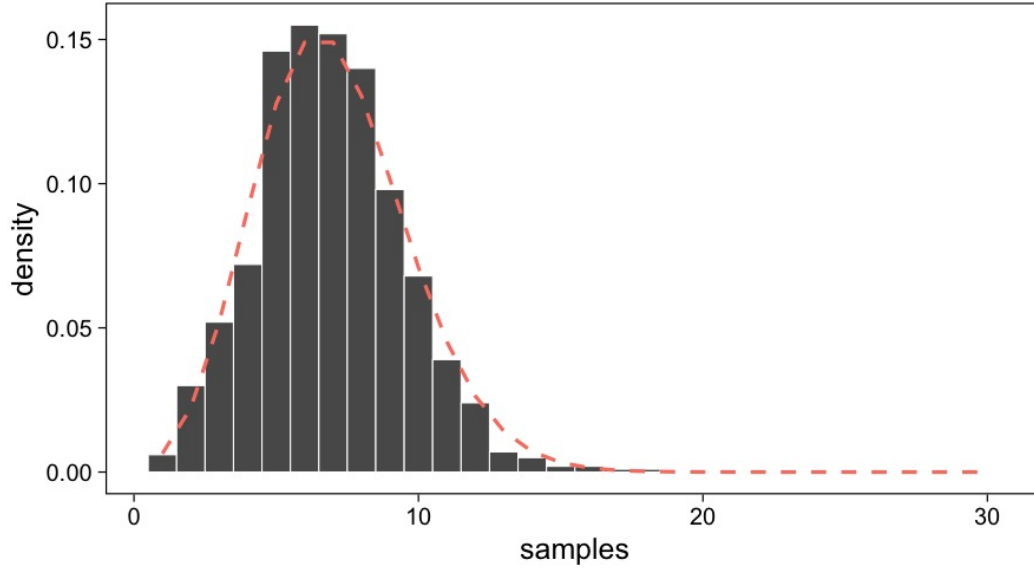
1 [1] 0.007295 0.029636 0.081765 0.172992 0.300708 0.449711 0.598714 0.729091
2 [9] 0.830496 0.901479 0.946650 0.973000 0.987189 0.994283 0.997593 0.999042
3 [17] 0.999638 0.999870 0.999956 0.999986 0.999995 0.999999 1.000000 1.000000

```

```

1 rpoisson <- function(nsamples, lambda){
2   probs <- ppois(1:30, lambda)
3   x <- c()
4   for (jj in 1:nsamples){
5     u <- runif(1)
6     x[jj] <- which(probs > u)[1]
7   }
8   return(x)
9 }

```



El problema de generar números aleatorios de la manera anterior es la necesidad de *guardar* el vector de probabilidades. Por ejemplo, una $\text{Poisson}(100)$. El intervalo $\lambda \pm 3\sqrt{\lambda}$ es $(70, 130)$.

3.2.1. Propiedad [Regla Empírica o regla de Pukelsheim [1]]: Si X es una variable aleatoria con media y varianza finitas. Entonces, la probabilidad de que una realización de X se encuentre a más de 3 desviaciones estándar de la media es a lo más 5 %.

4. MEZCLAS

Otra familia de distribuciones que es muy interesante de simular son las mezclas. Es decir, cuando podemos escribir

$$\pi(x) = \int_{\mathcal{Y}} \pi(x|y) \pi(y) dy, \quad \text{o} \quad \mathbb{P}(x) = \sum_{i \in \mathcal{Y}} \mathbb{P}(x|Y=i) \mathbb{P}(Y=i), \quad (10)$$

siempre y cuando sea sencillo generar números aleatorios de las marginales.

Por ejemplo, para generar números aleatorios de una t Student con ν grados de libertad. Podemos usar la representación

$$X|y \sim \mathcal{N}(0, \nu/y), \quad Y \sim \chi_{\nu}^2. \quad (11)$$

4.1. Binomial negativa

La variable aleatoria $X \sim \text{BinNeg}(n, \theta)$ tiene una representación

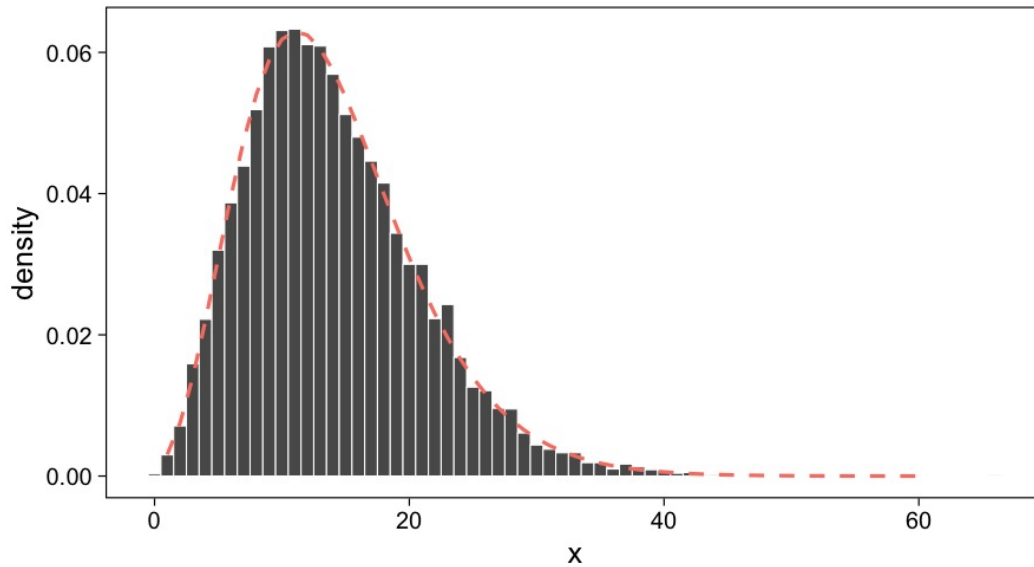
$$X|y \sim \text{Poisson}(y), \quad Y \sim \text{Gamma}(n, \beta), \quad (12)$$

donde $\beta = (1 - \theta)/\theta$.

```

1 nsamples <- 10^4
2 n <- 6; theta <- .3
3 y <- rgamma(nsamples, n, rate = theta/(1-theta))
4 x <- rpois(nsamples, y)

```



4.1.1. *Tarea* Prueba que una binomial negativa puede ser expresada como una mezcla Poisson-Gamma.

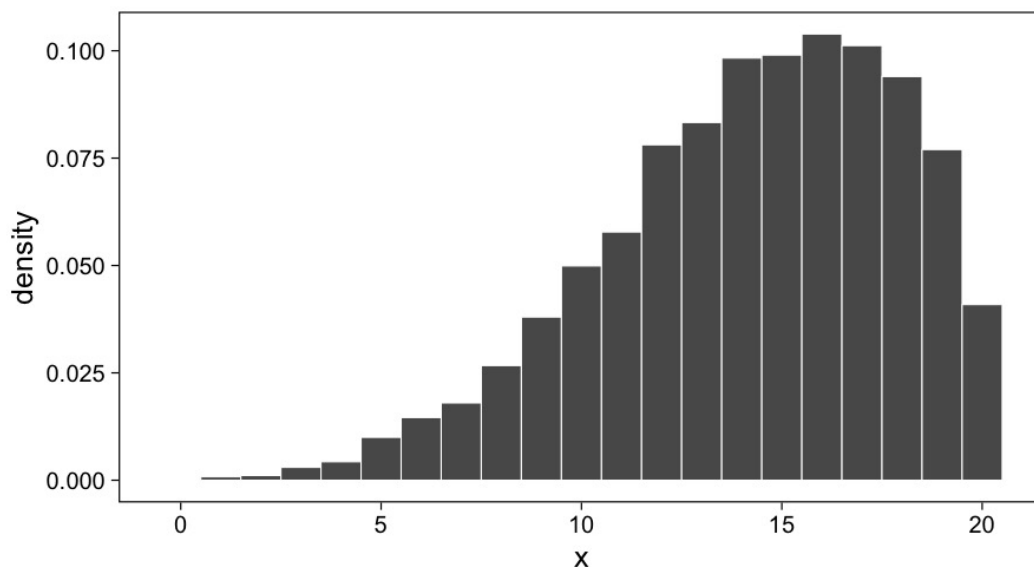
4.2. Beta Binomial

Una distribución bastante conocida es la distribución Beta-Binomial la cual, como su nombre indica, está conformada por una mezcla de una $\text{Bin}(n, \theta)$ y una $\text{Beta}(\alpha, \beta)$. Es decir, $x \sim \text{BetaBin}(n, \alpha, \beta)$ si su función de masa de probabilidad está dada por

$$\mathbb{P}(x) = \binom{n}{x} \frac{B(x + \alpha, n - x + \beta)}{B(\alpha, \beta)}, \quad (13)$$

donde $B(\alpha, \beta)$ se conoce como la función Beta $B(\alpha, \beta) = (\Gamma(\alpha)\Gamma(\beta))/\Gamma(\alpha + \beta)$.

```
1 nsamples <- 10^4
2 n <- 20; a <- 5; b <- 2
3 theta <- rbeta(nsamples, a, b)
4 x <- rbinom(nsamples, n, theta)
```



4.3. Mezclas Gaussianas

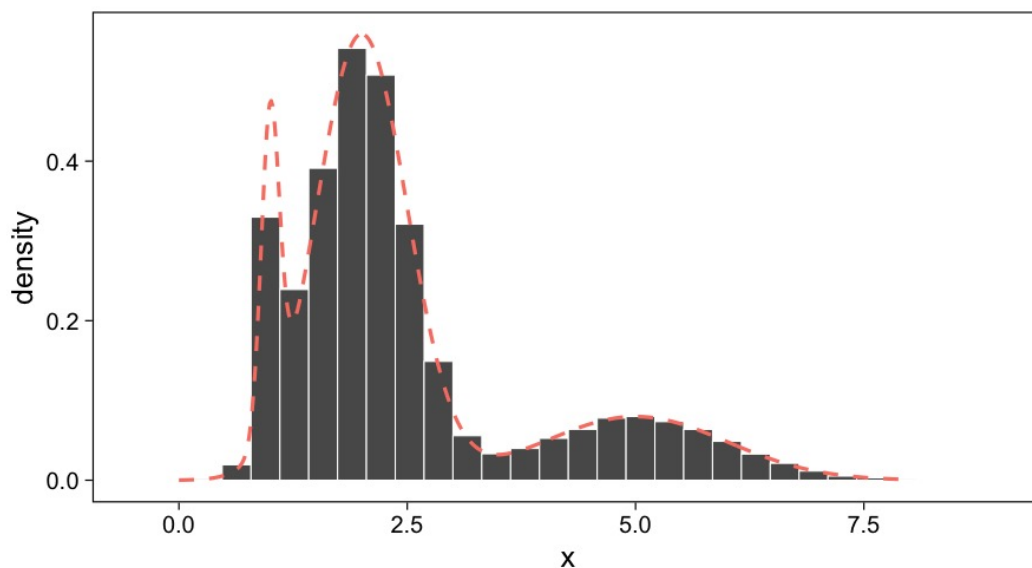
Otro modelo *famoso* es el de una mezcla de Gaussianas. Donde tenemos k posibles poblaciones cada una con proporción $\pi_k \in (0, 1)$ y $\sum \pi_k = 1$. Y además, cada población tiene comportamiento distinto

$$\pi(x|k) = N(x|\mu_k, \sigma_k). \quad (14)$$

```

1 nsamples <- 10^5
2 y <- sample(1:3, size = nsamples, prob = c(.1, .7, .2), replace = TRUE)
3 x <- rnorm(nsamples,
4           mean = ifelse(y==1, 1, ifelse(y==2, 2, 5)),
5           sd = ifelse(y==1, 0.1, ifelse(y==2, 0.5, 1)))

```



4.4. Comentarios

- El procedimiento para generar aleatorios de una mezcla nos permite obtener **distribuciones marginales a partir de la conjunta**.
- La combinación de componentes nos permite ser selectivos en los métodos de generación que podamos utilizar (respetando la estructura de **dependencia condicional**).
- Por ejemplo, en el modelo

$$\mathbb{P}(x) = \sum_{j \in \mathcal{Y}} \mathbb{P}(X|Y = j) \mathbb{P}(Y = j). \quad (15)$$

5. PRUEBA χ^2

Podemos usar otro mecanismo para probar estadísticamente si nuestros números pseudo aleatorios siguen la distribución que deseamos.

Podemos pensar en esta alternativa como la versión **discreta** de la prueba KS.

Lo que estamos poniendo a prueba es

$$H_0 : \mathbb{P}(x) = \mathbb{P}_0(x) \quad \forall x \quad \text{contra} \quad H_1 : \mathbb{P}(x) \neq \mathbb{P}_0(x) \quad \text{para alguna } x. \quad (16)$$

5.1.1. Procedimiento de la prueba χ^2

1. Hacemos una partición del rango de la distribución supuesta en k subintervalos con límites $\{a_0, a_1, \dots, a_k\}$, y definimos N_j como el número de observaciones (de nuestro generador de pseudo-aleatorios) en cada subintervalo.
2. Calculamos la proporción esperada de observaciones en el intervalo $(a_{j-1}, a_j]$ como

$$p_j = \int_{a_{j-1}}^{a_j} d\mathbb{P}(x). \quad (17)$$

3. La estadística de prueba es

$$\chi^2 = \sum_{j=1}^k \frac{(N_j - np_j)^2}{np_j}. \quad (18)$$

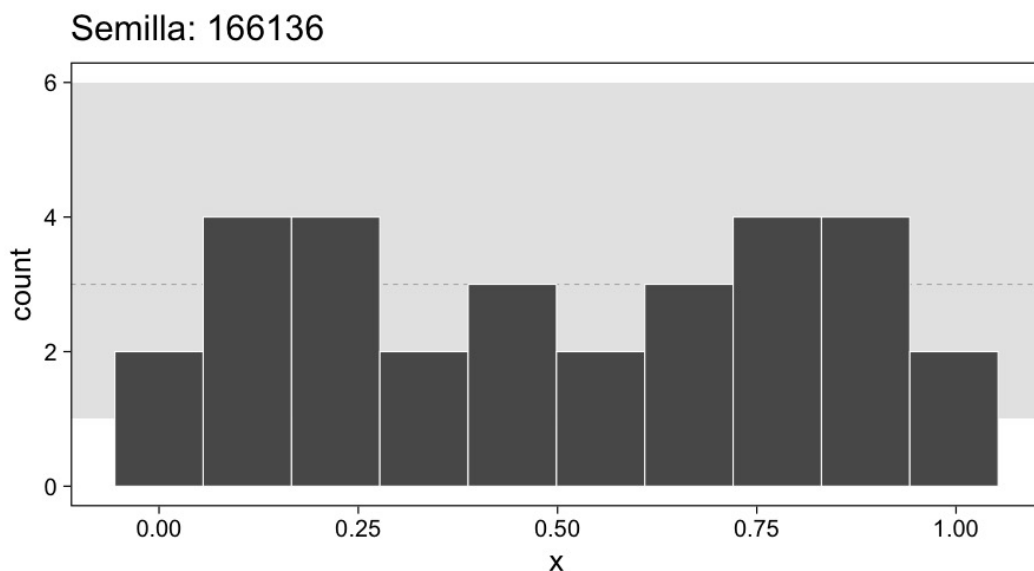
Nota que estamos comparando dos histogramas. El histograma observado que construimos a partir de nuestros números pseudo-aleatorios contra el histograma que esperaríamos de la distribución. ¿Puedes pensar en algún problema con esta prueba?

La visualización correspondiente sería lo siguiente. Utilizamos nuestro generador para obtener muestras.

```

1  ## Esto es para poner a prueba un pseudo generador
2  rpseudo.uniform <- function(nsamples, seed = 108727){
3    x0 <- seed; a <- 7**5; m <- (2**31)-1;
4    x <- x0;
5    for (jj in 2:nsamples){
6      x[jj] <- (a * x[jj-1]) %% m
7    }
8    x/m
9  }

```



5.1.1.1. Pregunta: ¿Qué esperaríamos de nuestro estadístico χ^2 si nuestro generador de pseudo-aleatorios es incorrecto?

5.2. Aplicación de la prueba

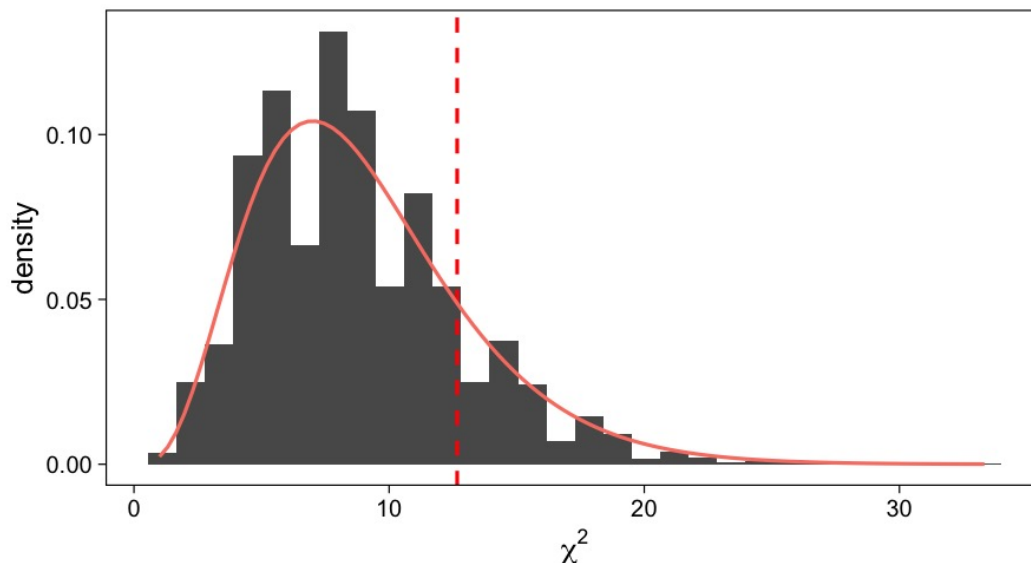
```

1 nsamples <- 30; nbreaks <- 10
2 samples <- data.frame(x = rpseudo.uniform(nsamples))
3
4 Fn <- hist(samples$x, breaks = nbreaks, plot = FALSE)$counts/nsamples
5 F0 <- 1/nbreaks
6
7 X2.obs <- (nsamples*nbreaks)*sum((Fn - F0)**2)

1 ## Esto es para generar datos observados de la distribucion que queremos
2 experiment <- function(nsamples){
3   nbreaks <- 10
4   samples <- data.frame(x = runif(nsamples))
5   Fn <- hist(samples$x, breaks = nbreaks, plot = FALSE)$counts/nsamples
6   F0 <- 1/nbreaks
7   X2 <- (nsamples*nbreaks)*sum((Fn - F0)**2)
8   return(X2)
9 }
10
11 X2 <- c()
12 for (jj in 1:5000){
13   X2[jj] <- experiment(nsamples)
14 }

```

En la Section 5.2 se muestra el histograma de las réplicas del estadístico χ^2 bajo el generador uniforme (lo tomamos como la distribución de la hipótesis nula) y comparamos contra el observado (línea punteada). Adicional, se incorpora la densidad de una χ^2_{k-1} (léase ji-cuadrada con $k - 1$ grados de libertad) que es la distribución asintótica del estadístico.



Por lo tanto, la probabilidad de haber observado un estadístico χ^2 tan extremo como el que observamos si el generador hubiera sido el que suponemos es:

```

1 [1] "Estadístico: 12.6667, Probabilidad: 0.1694"

```

Que podemos comparar contra el que obtenemos de una prueba "tradicional":

```

1 counts.obs ← Fn*nsamples
2 chisq.test(counts.obs, p = rep(1, nbreaks)/nbreaks, simulate.p.value = TRUE)

```

```

1
2      Chi-squared test for given probabilities with simulated p-value (based
3      on 2000 replicates)
4
5 data:  counts.obs
6 X-squared = 13, df = NA, p-value = 0.2

```

- La prueba χ^2 pues usualmente no es buena cuando el número de observaciones es menor a 50.
- La prueba KS tiene mejor potencia que la prueba χ^2 :

```

1 ks.test(samples$x, "punif")

```

```

1
2      Exact one-sample Kolmogorov-Smirnov test
3
4 data:  samples$x
5 D = 0.16, p-value = 0.4
6 alternative hypothesis: two-sided

```

5.3. Aplicación de pruebas

En la práctica se utiliza una colección de pruebas pues cada una es sensible a cierto tipo de desviaciones. La batería de pruebas mas utilizada es la colección de pruebas **DieHARD** que desarrolló [George Marsaglia](#) y que se ha ido complementando con los años.

6. DISTRIBUCIÓN NORMAL

Hasta ahora no hemos realmente hablado de cómo generar números de una de las distribuciones mas utilizadas. La distribución normal. Si pensamos en utilizar el método de la transformada inversa tendríamos que resolver la inversa de

$$\text{Prob}\{X \leq x\} = \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(s-\mu)^2}{2\sigma^2}} ds. \quad (19)$$

Sólo veremos ideas generales de cómo generamos números aleatorios de una distribución normal estándar.

6.1. Utilizando el teorema del límite central

Podemos utilizar una colección $U_1, \dots, U_k \stackrel{\text{iid}}{\sim} U(a, b)$ de tal forma que tengan

$$\mathbb{E}(U_i) = \frac{1}{2}, \quad \mathbb{V}(U_i) = \frac{1}{12}. \quad (20)$$

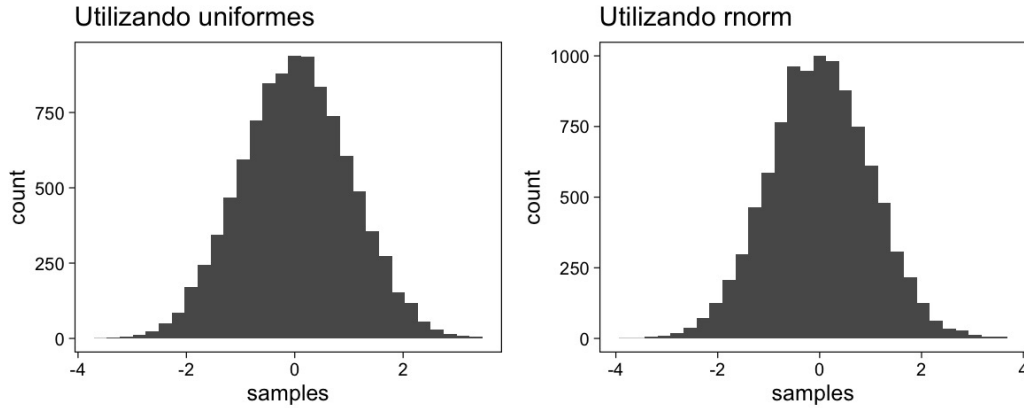
pues por el CLT tenemos que

$$Z = \frac{\sum_{i=1}^k U_j - k/2}{\sqrt{k/12}} \sim N(0, 1). \quad (21)$$

```

1 nsamples <- 10^4; k <- 12
2 U <- runif(k * nsamples)
3 U <- matrix(U, nrow = k)
4 X <- apply(U, 2, function(x){ (sum(x) - k/2)/sqrt(k/12) })

```



Depende de k la calidad de la aproximación. Sin embargo, siempre será un método aproximado.

6.2. El método de Box-Müller

Este método se basa en transformación de coordenadas polares a cartesianas. Es decir,

$$(u_1, u_2) \rightarrow (z_1, z_2), \quad (22)$$

donde u_i son uniformes independientes y z_i son normales estándar independientes.

Las coordenadas polares del vector (z_1, z_2) son

$$R^2 = z_1^2 + z_2^2, \quad \tan \theta = \frac{z_2}{z_1}. \quad (23)$$

La función de densidad conjunta para (R^2, θ) es igual a

$$\pi(r^2, \theta) = \frac{1}{2} \frac{1}{2\pi} e^{-\frac{r^2}{2}}, \quad 0 < r^2 < \infty, \quad 0 < \theta < 2\pi. \quad (24)$$

Podemos notar que tenemos

$$\pi(r^2, \theta) = \pi(r^2) \cdot \pi(\theta). \quad (25)$$

Es decir, son independientes.

Así que tenemos que

$$r^2 \sim \text{Exp}(2), \quad (26)$$

$$\theta \sim U(0, 2\pi). \quad (27)$$

Lo cual podemos finalmente escribir como

$$z_1 = \sqrt{-2 \log u_1} \cos(2\pi u_2), \quad (28)$$

$$z_2 = \sqrt{-2 \log u_1} \sin(2\pi u_2). \quad (29)$$

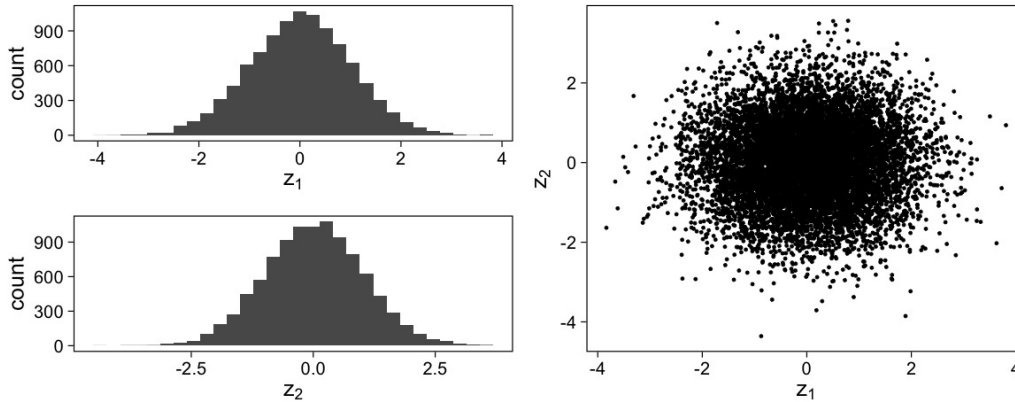
$$(30)$$

El problema de esto es que necesitamos calcular funciones complejas ($\sin(\cdot)$, $\cos(\cdot)$, $\sqrt{\cdot}$).

```

1 rnormal.bm ← function(n){
2   r ← sqrt(-2 * log(runif(n)))
3   theta ← runif(n, 0, 2 * pi)
4   z ← matrix(0, nrow = 2, ncol = n)
5   z[1,] ← r * cos(2 * pi * theta)
6   z[2,] ← r * sin(2 * pi * theta)
7   return(z)
8 }

```



6.3. Método de Marsaglia

Usaremos el método de Box-Müller como punto de partida (coordenadas polares y cartesianas). Lo que nos interesa es poder generar números aleatorios dentro del círculo. Así que generamos

$$V_1 = 2U_1 - 1, \quad V_2 = 2U_2 - 1, \quad (31)$$

a partir de $U_i \sim U(0, 1)$.

Si generamos un punto dentro del círculo unitario nos lo quedamos como una simulación válida. Si no, entonces repetimos la generación de los números aleatorios uniformes.

Si el punto se encuentra dentro del círculo entonces extraemos sus coordenadas por medio de

$$S = V_1^2 + V_2^2, \quad (32)$$

$$z_1 = \sqrt{\frac{-2 \log S}{S}} V_1, \quad z_2 = \sqrt{\frac{-2 \log S}{S}} V_2. \quad (33)$$

El método anterior tiene probabilidad $\pi/4$ de generar un punto dentro del círculo. Esto es, para generar 2 números gaussianos independientes usaremos en promedio $4/\pi \approx 1.273$ iteraciones de lanzar dardos. Esto es, generar en promedio 2.546 números aleatorios y un cálculo de logaritmos, una raíz cuadrada, una división y 4.546 multiplicaciones.

7. CONCLUSIONES Y COMENTARIOS

- El método de la transformada inversa es útil. Sin embargo, no siempre se puede calcular de forma cerrada (distribución normal).
- El modelo de mezclas es muy común en aplicaciones financieras, estudios de mercado y problemas de clasificación (segmentación).

REFERENCIAS

- [1] F. Pukelsheim. The three sigma rule. *The American Statistician*, 48(2):88–91, 1994. [6](#)
- [2] C. Robert and G. Casella. *Introducing Monte Carlo Methods with R*. Springer New York, New York, NY, 2010. [1](#)