

MAT4681 - Statistique pour les sciences

Arthur Charpentier

07 - Générer du hasard

été 2022

Générer du hasard I

A Million Random Digits with 100,000 Normal Deviates, 1955

★☆☆☆☆ Errors Throughout

Reviewed in the United States on December 30, 2013

They sure don't come up with random numbers like they used to. If you look closely, you will note that every tenth digit or so is just a repeat of the last digit and every hundredth or so is a just the same digit repeated three times. How sloppy!

A sampling of this "work":

Page 36 - Line 6 - 15 characters in should be 5, not 4.

Page 99 - Line 18 - first three characters should be "453" not "345".

Page 145 - Line 2 - 7th and 19th characters transposed.

Page 190 - Whole line of numbers omitted between 6th and 7th lines.

Pages 210 and 211 - Two sections appear quasi-randomized, instead of randomized.

★☆☆☆☆ Predictable

Reviewed in the United States on October 30, 2013

It seemed like about 10% of the time I was able to predict which number was next. It was still better than Life of Pi which, aside from being irrational, included no estimations of Pi at all.

★☆☆☆☆ Not really random

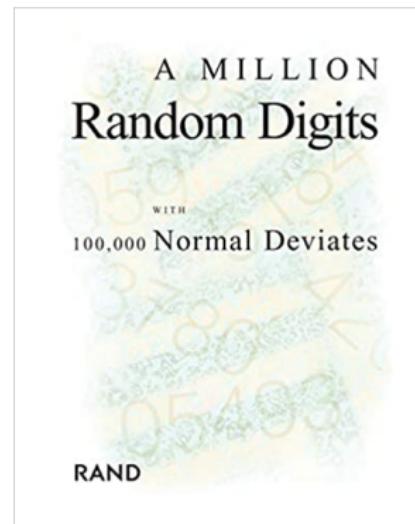
Reviewed in the United States on September 26, 2012

I bought two copies of this book. I find that the first copy perfectly predicts what the numbers will be in the second copy. I feel cheated.

★☆☆☆☆ great yet lacking....

Reviewed in the United States on May 1, 2013

So far, it's great, but what I REALLY want is an audio version that I can listen to at the beach or the gym. Maybe narrated by Morgan Freeman???



Générer du hasard II

A Million Random Digits with 100,000 Normal Deviates, 1955

★☆☆☆☆ Errors Throughout

Reviewed in the United States on December 30, 2013

They sure don't come up with random numbers like they used to. If you look closely, you will note that every tenth digit or so is just a repeat of the last digit and every hundredth or so is a just the same digit repeated three times.
How sloppy!

A sampling of this "work":

Page 36 - Line 6 - 15 characters in should be 5, not 4.

Page 99 - Line 18 - first three characters should be "453" not "345".

Page 145 - Line 2 - 7th and 19th characters transposed.

Page 190 - Whole line of numbers omitted between 6th and 7th lines.

Pages 210 and 211 - Two sections appear quasi-randomized, instead of randomized.

★☆☆☆☆ Predictable

Reviewed in the United States on October 30, 2013

It seemed like about 10% of the time I was able to predict which number was next. It was still better than Life of Pi which, aside from being irrational, included no estimations of Pi at all.

★☆☆☆☆ Not really random

Reviewed in the United States on September 26, 2012

I bought two copies of this book. I find that the first copy perfectly predicts what the numbers will be in the second copy. I feel cheated.

★☆☆☆☆ great yet lacking....

Reviewed in the United States on May 1, 2013

So far, it's great, but what I REALLY want is an audio version that I can listen to at the beach or the gym. Maybe narrated by Morgan Freeman???

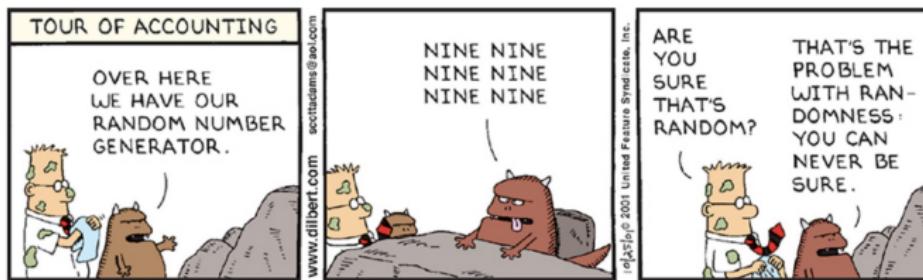
TABLE OF RANDOM 2141 TS

3

00100	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00101	38055 03354	32884 09780	08350 00860	09720 47763	71209 23603
00102	17945 72705	83018 04812	02017 03011	03710 07770	07709 00039
00103	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00104	60972 68777	39810 03956	34040 40819	20549 06041	33484 00479
00105	34126 60051	37699 06494	34485 46672	01958 77100	09899 70754
00106	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00107	30532 33704	10974 12202	30985 23309	00861 08882	53986 66480
00108	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00109	48258 45239	09790 47619	33135 07453	35660 33973	16818 06311
00110	40065 04852	38670 31348	42014 28297	05191 28351	08963 79231
00111	03799 42052	18682 04463	34194 41374	75071 14724	09958 18605
00112	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00113	19032 033845	17620 02862	05647 08846	78138 06850	19440 09412
00114	11220 94747	07210 37640	65309 22909	27482 45476	05144 05159
00115	31151 07782	08040 00219	13470 04359	23880 00001	00183 19121
00116	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00117	20204 47764	07681 48352	32785 06032	38097 39405	04236 01340
00118	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00119	78612 16929	33902 91236	44272 18140	05741 16014	47819 37220
00120	41489 04747	04800 10233	24977 05200	74810 04740	19325 81549
00121	46252 02049	02818 52466	45526 07531	61219 31180	14413 70951
00122	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00123	52701 08337	04830 07713	16330 09870	11864 09610	92181 56181
00124	57778 30696	04830 04869	46902 27776	02122 00000	02122 00000
00125	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00126	33683 84924	00410 09814	23167 08663	55882 47941	06222 45391
00127	02694 48297	39986 02115	53509 49967	46821 41375	48977 04627
00128	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00129	38608 32486	45134 03834	39004 72689	43917 51860	43652 56993
00130	20183 01069	70024 10921	41190 07231	71367 13087	08977 11403
00131	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00132	26815 43620	37765 05844	05844 00000	27703 14767	01477 02097
00133	64397 11062	03237 02142	20547 01759	45197 28312	03748 03967
00134	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00135	62761 00972	42393 01416	00832 20882	01291 47349	03948 07126
00136	14297 00420	06864 00271	42510 06218	05894 74791	41196 77460
00137	51329 92246	08988 77074	88722 36706	06168 49423	09912 03179
00138	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00139	05446 35096	03128 10846	74457 06061	72848 11834	79982 68616
00140	39528 72884	03474 25953	48045 33347	18619 13674	18611 19241
00141	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00142	07986 16120	02841 22820	21094 13141	33292 19763	01189 67940
00143	90795 04235	13930 02709	03986 04866	48265 00000	02147 02112
00144	00801 10145	03716 16894	06683 04852	08469 05323	08618 19181
00145	01704 20503	04737 21031	70561 00529	47086 43823	09912 03476
00146	34414 82137	06887 50567	19153 00023	12302 00783	32814 68091
00147	62429 75163	44989 16822	06004 00607	76379 41600	45981 74486
00148	78903 04748	52162 00286	41188 00000	49772 00000	02147 02112
00149	01704 20503	04737 21031	70561 00529	47086 43823	09912 03476

Générer du hasard III

voir aussi Philipps (2020)



(dessin via <https://dilbert.com/strip/2001-10-25>)

“A random sequence is a vague notion... in which each term is unpredictable to the uninitiated and whose digits pass a certain number of tests traditional with statisticians...” Derrick Lehmer, cité par Knuth (1997) (voir aussi L'Ecuyer (2017) pour une perspective historique)

Générer du hasard IV

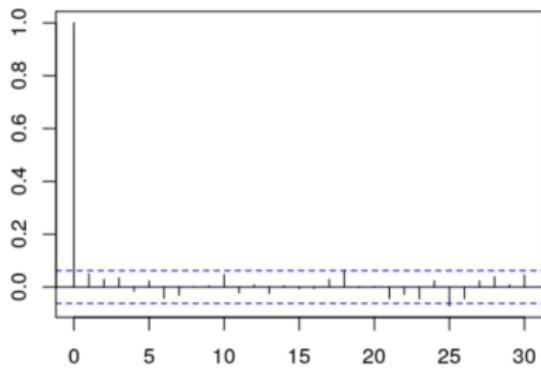
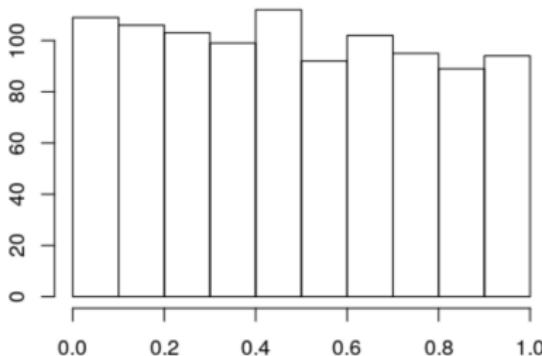
Heuristiquement, une suite (u_i) semble **aléatoire** si

1. les tirages sont **uniformes** sur $[0, 1]$: $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{u_i \in (a, b)} = b - a$

avec $b > a$,

2. les appels sont **indépendants**: pour $b > a$ et $d > c$.

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{u_i \in (a, b), u_{i+k} \in (c, d)} = (b - a)(d - c), \quad \forall k \in \mathbb{N}$$



Générer du hasard V

- ▶ approche physique

Utiliser des tirages de pièces, de dés, de roulette, etc

$$\text{PFPFFPPPF} \rightarrow 2^9 \cdot 1 + 2^8 \cdot 0 + 2^7 \cdot 1 + 2^6 \cdot 1 + 2^5 \cdot 0 + 2^4 \cdot 1 + 2^3 \cdot 1 + 2^2 \cdot 1 + 2^1 \cdot 0 + 2^0 \cdot 0 = 732$$

$$\text{PPPPPPFPPF} \rightarrow 2^9 \cdot 1 + 2^8 \cdot 1 + 2^7 \cdot 1 + 2^6 \cdot 1 + 2^5 \cdot 1 + 2^4 \cdot 1 + 2^3 \cdot 0 + 2^2 \cdot 1 + 2^1 \cdot 1 + 2^0 \cdot 0 = 1014$$

- ▶ Utiliser des techniques mathématiques

von Neumann (1949), “middle-square method”,

$$\begin{aligned} 8653^2 &= 74874409 \rightarrow 8744^2 = 76457536 \rightarrow 4575^2 = 20930625 \\ 9306^2 &= 86601636 \rightarrow 6016^2 = 36192256 \rightarrow 1922^2 = 3694084 \\ 6940^2 &= 48163600 \rightarrow 1636^2 = 2676496 \rightarrow 6764^2 = 45751696 \end{aligned}$$

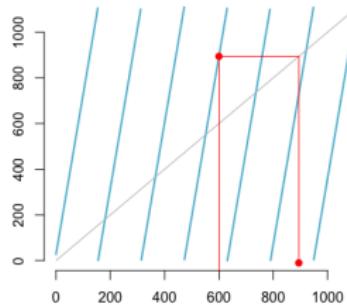
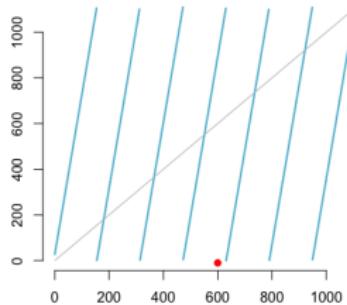
or use some linear congruential generator, from Lehmer (1951)

Générer du hasard VI

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = f(x_{n-1}) = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_0 = 600$, so that $x_1 = 894$

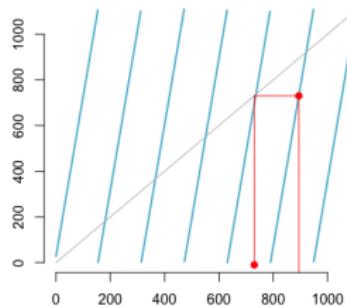
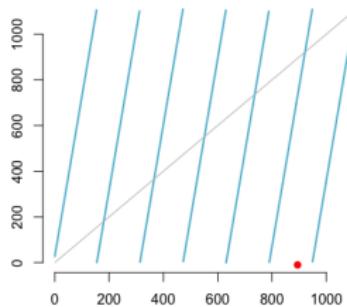


Générer du hasard VII

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_1 = 894$, so that $x_2 = 730$

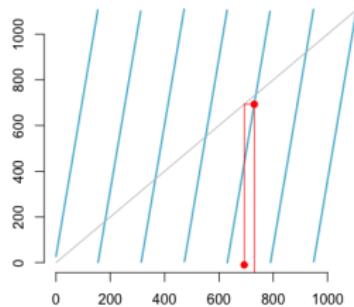
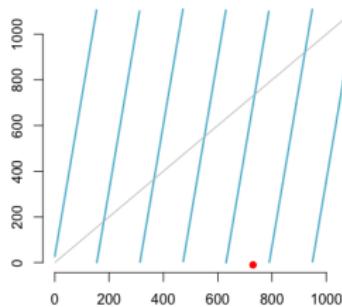


Générer du hasard VIII

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_2 = 730$, so that $x_3 = 693$

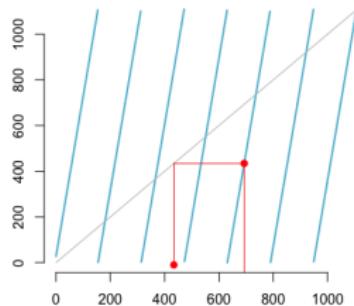
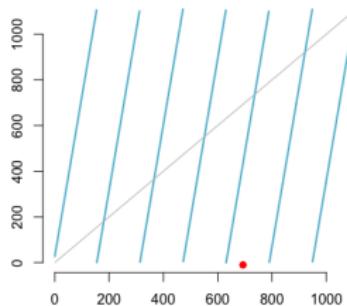


Générer du hasard IX

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_3 = 693$, so that $x_4 = 434$

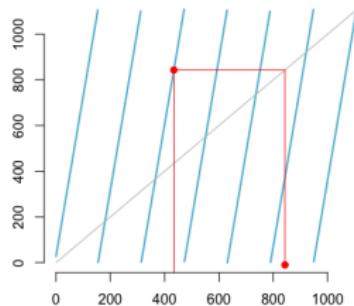
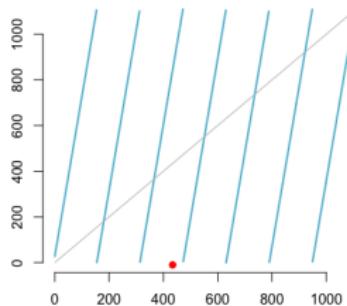


Générer du hasard X

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_4 = 434$, so that $x_5 = 843$

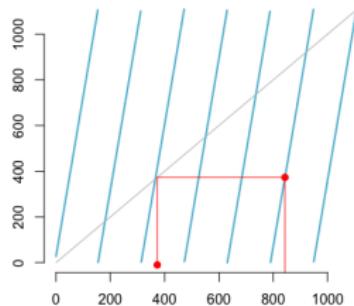
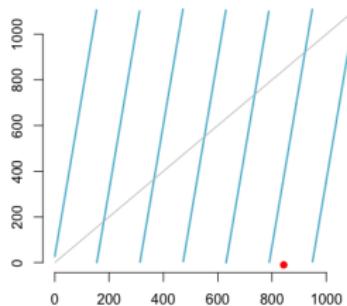


Générer du hasard XI

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_5 = 843$, so that $x_6 = 373$

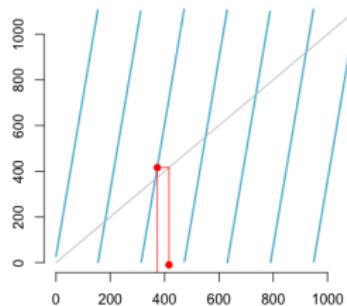
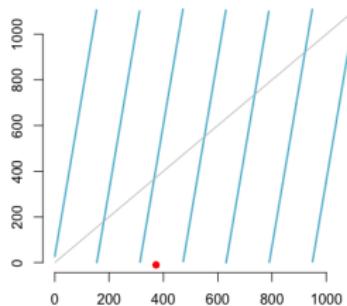


Générer du hasard XII

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_6 = 373$, so that $x_7 = 416$

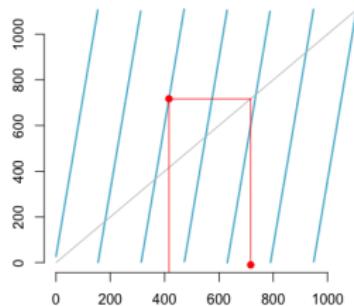
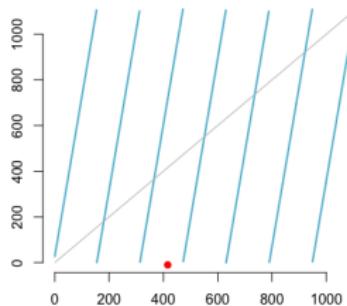


Générer du hasard XIII

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_7 = 416$, so that $x_8 = 717$

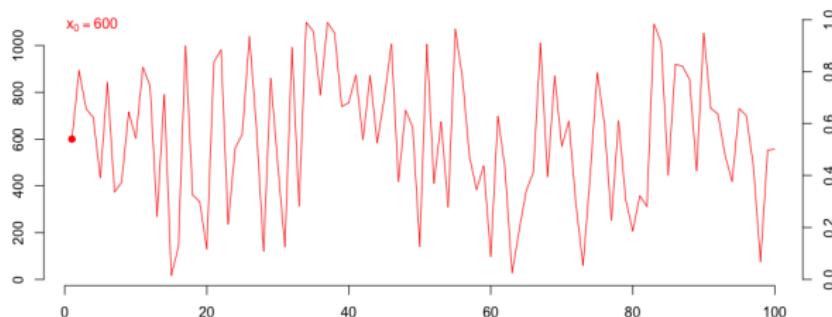


Générer du hasard XIV

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_0 = 600$, see x_1, x_2, \dots, x_{100}

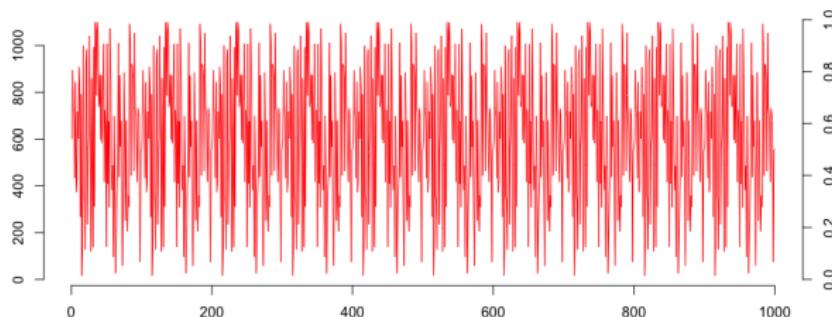


Générer du hasard XV

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_0 = 600$, see $x_1, x_2, \dots, x_{1000}$

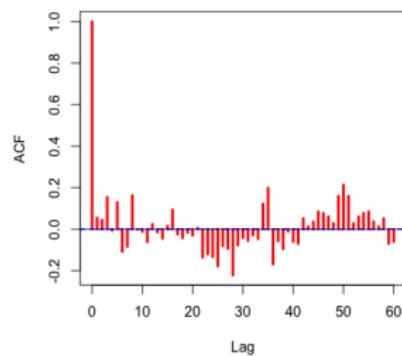
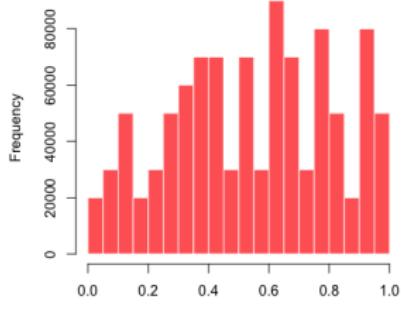


Générer du hasard XVI

Sequences generated with Sedgewick algorithm, $u_n = x_n / m$,

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, with $x_0 = 600$, distribution of x_i 's and autocorrelation

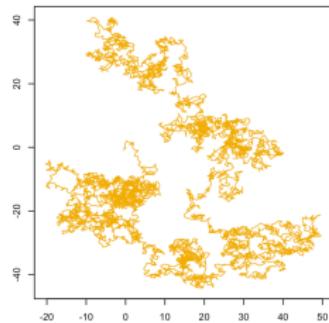
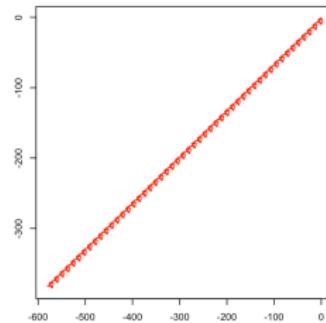
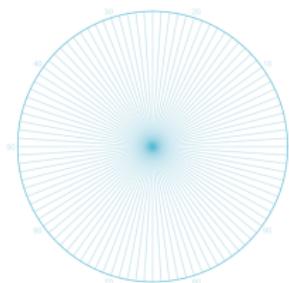


Générer du hasard XVII

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 7$, $c = 27$ and $m = 1111$, or $a = 74$, $c = 75$ and
 $m = 2^{16} + 1 = 65537$ (see ZX81),

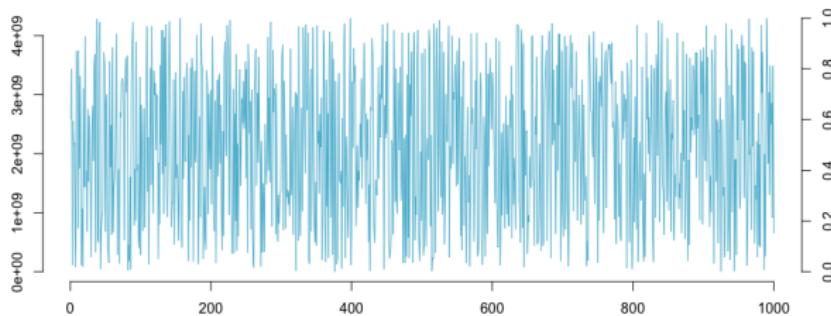


Générer du hasard XVIII

Sequences generated with Sedgewick algorithm, $u_n = x_n / m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 1013904223$, $c = 1664525$ and $m = 2^{32}$, with $x_0 = 2576980378$,

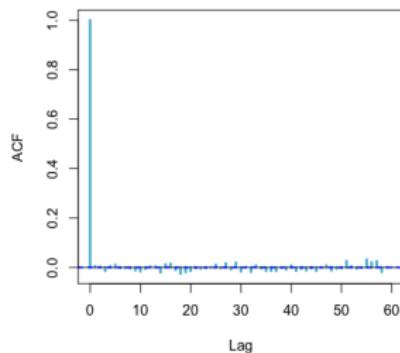
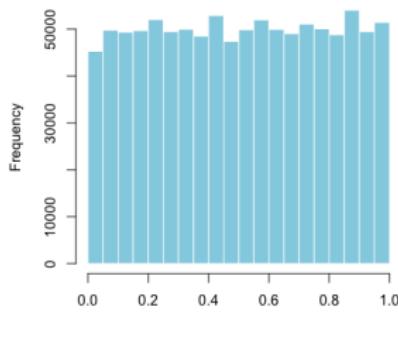


Générer du hasard XIX

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 1013904223$, $c = 1664525$ and $m = 2^{32}$, with $x_0 = 2576980378$,

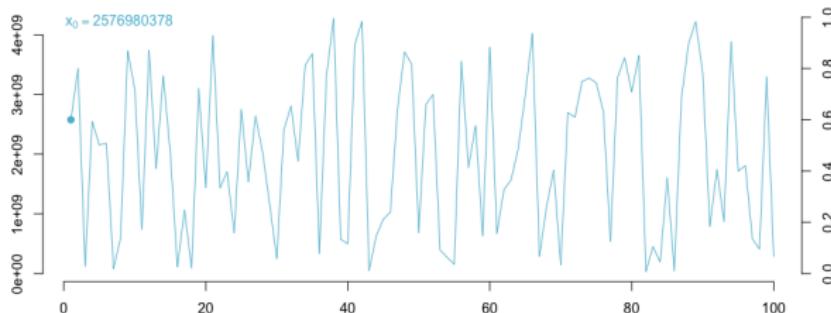


Générer du hasard XX

Sequences generated with Sedgewick algorithm, $u_n = x_n / m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 1013904223$, $c = 1664525$ and $m = 2^{32}$, with $x_0 = 2576980378$

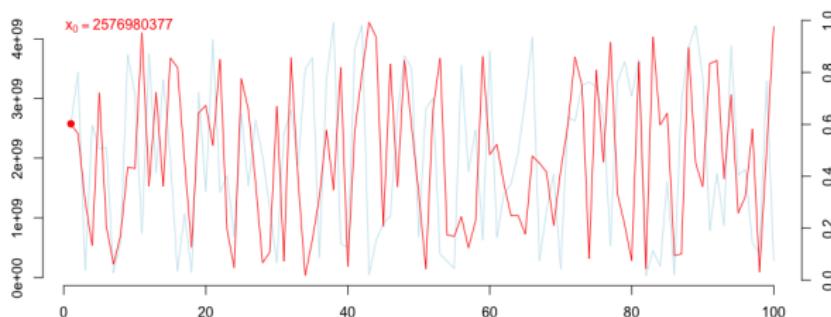


Générer du hasard XXI

Sequences generated with Sedgewick algorithm, $u_n = x_n / m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 1013904223$, $c = 1664525$ and $m = 2^{32}$, with $x_0 = 2576980377$

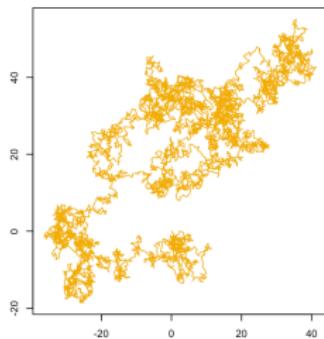
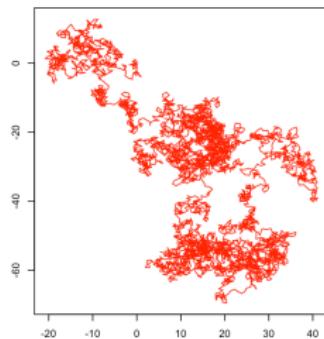
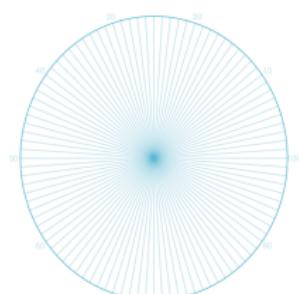


Générer du hasard XXII

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 1013904223$, $c = 1664525$ and $m = 2^{32}$,

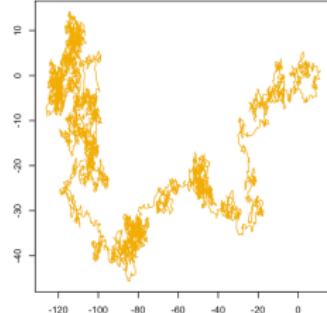
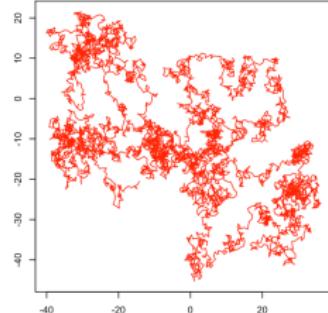
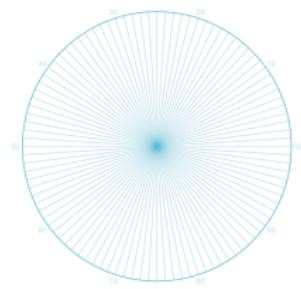


Générer du hasard XXIII

See also the use of π 's digits

3.14159265358979323846264338327950288419716939937510582097494459230
8164062862089986280348253421170679821480865132823066470938446095505

3.14159265358979323846264338327950288419716939937510582097494459230
8164062862089986280348253421170679821480865132823066470938446095505



Générer du hasard XXIV

A lot of connexions with **chaos theory**

"when the present determines the future, but the approximate present does not approximately determine the future" Edward Lorenz, cited in **May (1976)**

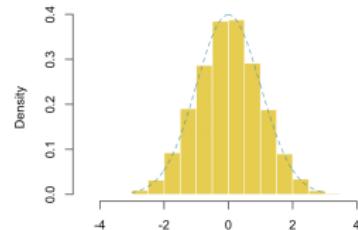
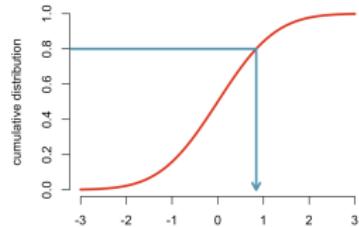
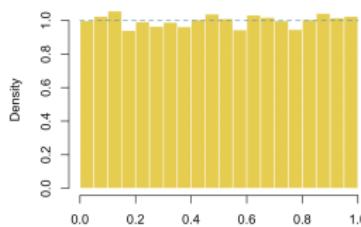
Monte Carlo I

Probability Inverse Transform

Let $F : \mathbb{R} \rightarrow [0, 1]$ denote a cumulative distribution function. If $U \sim \mathcal{U}([0, 1])$, then $X = F^{-1}(U)$ has cumulative distribution function F , where $F^{-1}(u) = \inf\{x : F(x) \geq u\}$.

Proof: Let $x \in \mathbb{R}$, $\mathbb{P}[X \leq x]$ is equal to

$$\mathbb{P}[F^{-1}(U) \leq x] = \mathbb{P}[F(F^{-1}(U)) \leq F(x)] = \mathbb{P}[U \leq F(x)] = F(x)$$



Monte Carlo II

Law of large numbers - Central limit theorem

$$\frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n \xrightarrow{\mathbb{P}} \mu = \mathbb{E}_{\mathbb{P}}[X], \text{ i.e. } \forall \varepsilon, \lim_{n \rightarrow \infty} \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| < \varepsilon\right] = 1$$

$$\frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n \xrightarrow{\text{a.s.}} \mu = \mathbb{E}_{\mathbb{P}}[X], \text{ i.e. } \forall \varepsilon, \mathbb{P}\left[\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right] = 1$$

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma^2) \text{ where } \mu = \mathbb{E}_{\mathbb{P}}[X] \text{ and } \sigma^2 = \mathbb{E}_{\mathbb{P}}[(X - \mu)^2]$$

Interestingly, the central limit theorem is valid also in higher dimension

$$\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \boldsymbol{\Sigma}) \text{ where } \boldsymbol{\mu} = \mathbb{E}_{\mathbb{P}}[\mathbf{X}] \text{ and } \boldsymbol{\Sigma} = \mathbb{E}_{\mathbb{P}}[(\mathbf{X} - \boldsymbol{\mu}) \times (\mathbf{X} - \boldsymbol{\mu})]$$

(speed of convergence does not depend on the dimension)

Monte Carlo III

If (U_i) are uniform independent random variables, if $g : \mathbb{R} \rightarrow \mathbb{R}$,

$$\frac{1}{n} \sum_{i=1}^n g(U_i) \xrightarrow{\mathbb{P}} \mathbb{E}_{\mathbb{P}}[g(U)] = \int_0^1 g(u) d\mathbf{u}, \text{ where } U \sim \mathcal{U}([0, 1]).$$

If (\mathbf{U}_i) are uniform independent random vectors, with

$\mathbf{U}_i = (U_{i,1}, \dots, U_{i,d})$, if $g : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\frac{1}{n} \sum_{i=1}^n g(\mathbf{U}_i) \xrightarrow{\mathbb{P}} \mathbb{E}_{\mathbb{P}}[g(\mathbf{U})] = \int_{[0,1]^d} g(\mathbf{u}) d\mathbf{u}, \text{ where } \mathbf{U} \sim \mathcal{U}([0, 1]^d).$$

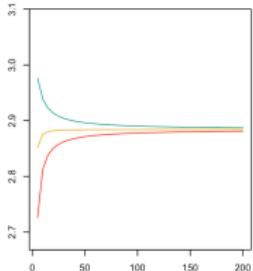
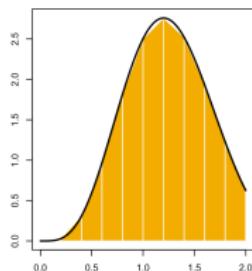
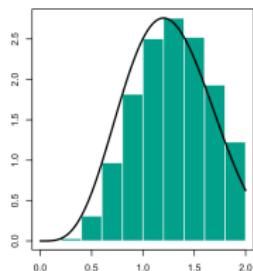
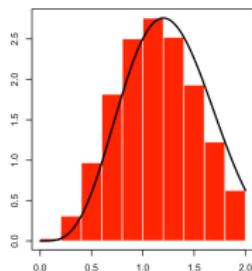
Note that either in dimension 1 or in dimension d , the speed of convergence is $O(n^{-1/2})$.

Monte Carlo IV

Example: integral approximation, error bounded in $O(n^{-1})$,
 $O(n^{-1})$ or $O(n^{-2})$

$$\int_0^2 h(y)dy = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{2}{n} \times h\left((i-1)\frac{2}{n}\right) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{2}{n} \times h\left(i\frac{2}{n}\right)$$

$$\int_0^2 h(y)dy = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{2}{n} \times \frac{1}{2} \left(h\left((i-1)\frac{2}{n}\right) + h\left(i\frac{2}{n}\right) \right)$$



Monte Carlo V

$$\left| \int_0^2 h(y) dy - \sum_{i=1}^n \frac{2}{n} \times h\left(i \frac{2}{n}\right) \right| \leq \frac{2^2}{n} \|f'\|_\infty$$

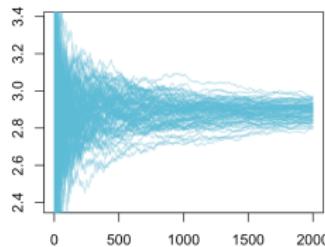
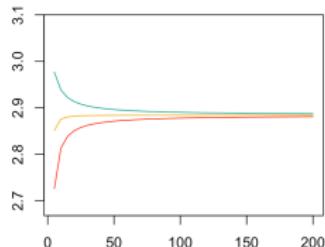
Using the law of large numbers

$$\int_0^2 h(y) dy = \int_0^1 2h(2u) du \lim_{n \rightarrow \infty} \frac{2}{n} \sum_{i=1}^n h(2u_i)$$

Convergence is slow,

$$\mathbb{P} \left[\left| \int_0^2 h(y) dy - \frac{2}{n} \sum_{i=1}^n h(2u_i) \right| \leq \frac{1.96\sigma}{\sqrt{n}} \right] = 95\%$$

which is $O(n^{-1/2})$, with a given probability.



How to create randomness?

Linear Congruential Method

Given $a, b, m \in \mathbb{N}_*$ and $x_0 \in \{0, 1, \dots, m\}$, define

$$x_{i+1} = (ax_i + b) \text{ modulo } m,$$

and set $u_i = x_i / m$.

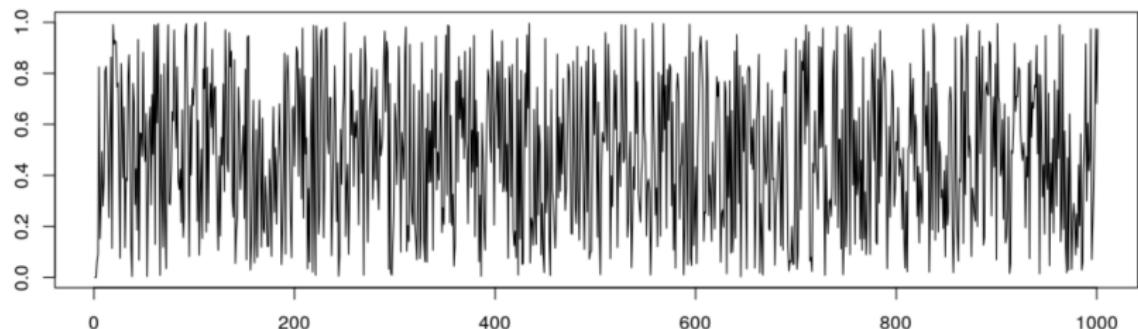
```
1 > a = 13; b = 43; m = 100; x = 77; u = rep(NA, 40)
2 > for (i in 1:40) {x = (a * x + b) %% m
3 +     u[i] = x / m }
4 > u
5 [1] 0.44 0.15 0.38 0.37 0.24 0.55 0.58 0.97 0.04 0.95
6 [11] 0.78 0.57 0.84 0.35 0.98 0.17 0.64 0.75 0.18 0.77
7 [21] 0.44 0.15 0.38 0.37 0.24 0.55 0.58 0.97 0.04 0.95
```

Problem: not all values in $\{0, \dots, m-1\}$ are obtained, and there is a cycle here.

How to create randomness?

Solution: (very) large values for m and choose properly a and b .

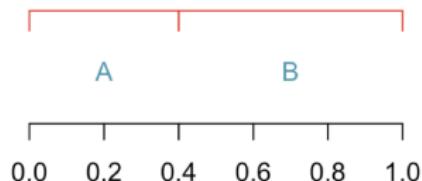
E.g. $m = 2^{32} - 1$, $a = 16807$ ($= 7^5$) and $b = 0$ (used in Matlab).



Génération de loi binomiale

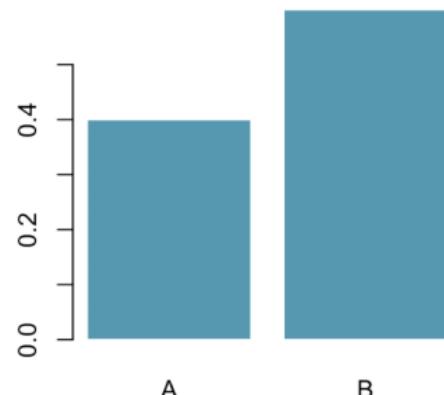
Soit $p \in (0, 1)$, $U \sim \mathcal{U}_{[0,1]}$

$$X = \begin{cases} 1 & \text{si } U < p \\ 0 & \text{si } U \geq p \end{cases}$$



```
1 > p = 0.4
2 > n = 1e7
3 > U1 = runif(n)
4 > Z = (U1< p)*1
5 > barplot(table(Z)/n)
```

```
1 > Z = sample(0:1, size=n,
   replace=TRUE ,
2   prob=c(1-p,p))
3 > table(Z)
Z
5      0          1
6 6000144 3999856
```



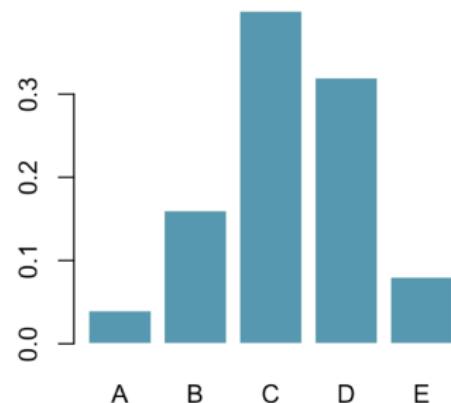
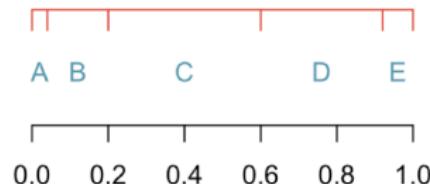
Génération de loi multinomiale

Soit \mathbf{p} un vecteur de probabilité,

$$\bar{p}_1 = 0 \text{ et } \bar{p}_{j+1} = \sum_{i=1}^j p_i$$

$$U \in [\bar{p}_j, \bar{p}_{j+1}) \implies X = j + 1$$

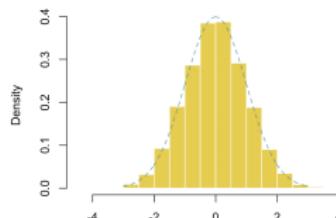
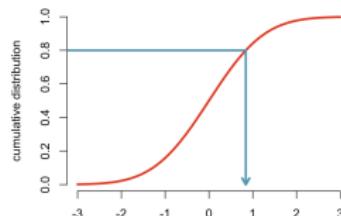
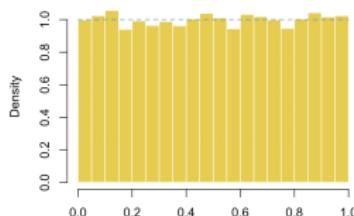
```
1 > p = c(0.04, 0.16, 0.40,
      0.32, 0.08)
2 > cumsum(p)
3 [1] 0.04 0.20 0.60 0.92 1.00
4 > n = 1e7
5 > U1 = runif()
6 > Z = cut(U1, breaks = c(0,
      cumsum(p)), labels =
      LETTERS[1:5])
7 > barplot(table(Z)/n)
```



Inversion de la Fonction de Répartition

Inversion de fonction de répartition, inverse method sampling

Soit F une fonction de répartition, si $U \sim \mathcal{U}([0, 1])$, $X = F^{-1}(U)$ a pour fonction de répartition F .



F^{-1} est simplement la fonction quantile,

```
1 > Q = function(u) qnorm(u, 0, 1)
```

Inversion de la Fonction de Répartition

```
1 > U = runif(100)
2 [1] 0.27 0.37 0.57 0.91 0.20 0.90 0.94 0.66 0.63 0.06
3 [11] 0.21 0.18 0.69 0.38 0.77 0.50 0.72 0.99 0.38 0.78
4 [21] 0.93 0.21 0.65 0.13 0.27 0.39 0.01 0.38 0.87 0.34
5 [31] 0.48 0.60 0.49 0.19 0.83 0.67 0.79 0.11 0.72 0.41
6 [41] 0.82 0.65 0.78 0.55 0.53 0.79 0.02 0.48 0.73 0.69
7 [51] 0.48 0.86 0.44 0.24 0.07 0.10 0.32 0.52 0.66 0.41
8 [61] 0.91 0.29 0.46 0.33 0.65 0.26 0.48 0.77 0.08 0.88
9 [71] 0.34 0.84 0.35 0.33 0.48 0.89 0.86 0.39 0.78 0.96
```

```
1 > Q(U)
2 [1] -0.63 -0.33 0.18 1.33 -0.84 1.27 1.60 0.41
3 [9] 0.33 -1.54 -0.82 -0.93 0.49 -0.29 0.74 -0.01
4 [17] 0.58 2.40 -0.31 0.76 1.51 -0.80 0.39 -1.15
5 [25] -0.62 -0.29 -2.21 -0.30 1.12 -0.41 -0.04 0.25
6 [33] -0.02 -0.89 0.94 0.44 0.82 -1.24 0.59 -0.22
7 [41] 0.92 0.38 0.78 0.13 0.07 0.80 -1.99 -0.06
8 [49] 0.62 0.50 -0.06 1.09 -0.16 -0.69 -1.47 -1.28
9 [57] -0.48 0.05 0.42 -0.24 1.36 -0.54 -0.10 -0.43
10 [65] 0.39 -0.65 -0.05 0.73 -1.38 1.15 -0.41 0.99
```

Inversion de la Fonction de Répartition Empirique

Given a sample $\{x_1, \dots, x_n\}$ i.i.d. from F ,

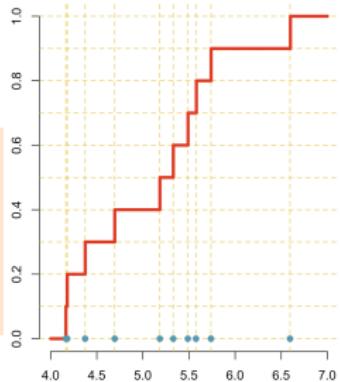
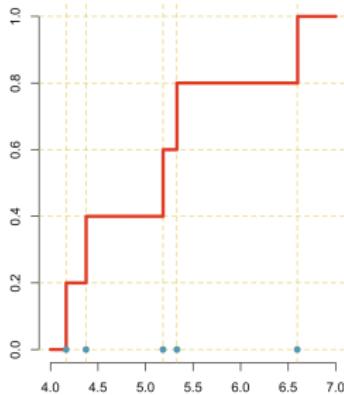
$$F(x) = \mathbb{P}[X \leq x],$$

the empirical cumulative distribution function is

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(x_i \leq x), \quad x \in \mathbb{R}$$

Glivenko-Cantelli: $\hat{F}_n \rightarrow F$ as $n \rightarrow \infty$.

```
1 > Finvemp = function(u,x) sort(x)[  
2   ceiling(u*length(x))]  
3 > Qinv = Vectorize(function(u){  
4   Finvemp(u,x)  
5 })
```



Inversion de la Fonction de Répartition Empirique

The inverse method with \hat{F}_n simply means resampling within $\{x_1, \dots, x_n\}$ with equal probabilities $1/n$ (or *with replacement*)

```
1 > x
2 [1] 4.164 4.374 5.184 5.330 6.595
3 > Qemp(U)
4 [1] 6.60 6.60 6.60 5.33 4.37 5.33 5.33 4.16 6.60 5.33
5 [11] 4.37 4.37 4.37 6.60 5.33 5.18 5.33 5.18 6.60 5.18
6 [21] 5.18 4.37 6.60 4.37 4.16 6.60 4.16 6.60 5.33 4.16
7 [31] 4.16 6.60 4.37 4.37 5.33 5.18 5.18 5.18 5.33 5.33
8 [41] 4.37 5.18 5.33 5.18 4.37 5.18 5.18 5.18 5.33 5.18
9 [51] 5.33 4.37 4.37 4.16 5.18 5.18 5.18 5.18 4.16 5.18
10 [61] 4.37 4.16 4.16 4.16 6.60 4.37 4.37 5.33 5.18 4.16
11 [71] 5.33 4.16 6.60 5.18 4.16 4.16 5.18 4.16 5.18 4.16
```

called **bootstrapping**

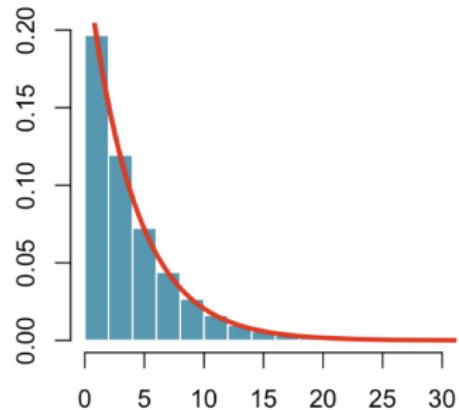
```
1 > sample(x, size = 80, replace = TRUE)
```

Génération de loi Exponentielle

$F(x) = \mathbb{P}[X < x] = 1 - e^{-ax}$ pour $x \geq 0$. On veut $1 - e^{-aq} = u$,
i.e. $e^{-aq} = 1 - u$, $aq = -\log(1 - u)$

$$F^{-1}(u) = \frac{-1}{a} \log(1 - u), \text{ pour } u \in [0, 1].$$

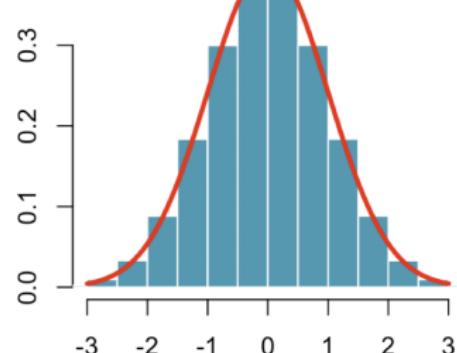
```
1 > a = 1/4
2 > U1 = runif(1e7)
3 > Z = -log(1-U1)/a
4 > hist(Z, probability=TRUE,
      xlim=c(0,30))
5 > curve(dexp(x,a), add=TRUE)
```



Génération de loi Gaussienne $\mathcal{N}(0, 1)$

Si $U_1, U_2 \sim \mathcal{U}_{[0,1]}$, indépendantes, $R = \sqrt{-2 \log(U_1)}$ et $\Theta = 2\pi U_2$, alors $(X_1, X_2) = (R \cos \Theta, R \sin \Theta)$ est un couple de variables $\mathcal{N}(0, 1)$ indépendantes

```
1 > U1 = runif(1e7)
2 > U2 = runif(1e7)
3 > R = sqrt(-2*log(U1))
4 > Theta = 2*pi*U2
5 > Z = R*cos(Theta)
6 > hist(Z, proba=TRUE)
7 > curve(dnorm(x, 0, 1), add=TRUE)
```

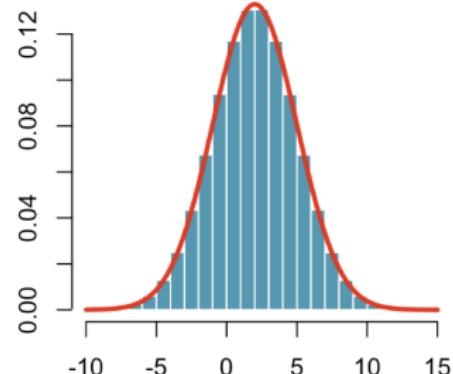


On parle de la méthode de Box-Muller.

Génération de loi Gaussienne $\mathcal{N}(\mu, \sigma^2)$

Si $Z \sim \mathcal{N}(0, 1)$, $X = \mu + \sigma Z \sim \mathcal{N}(\mu, \sigma^2)$

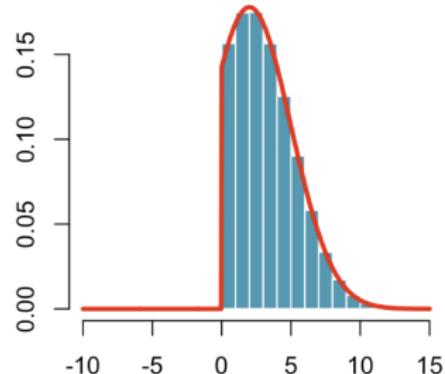
```
1 > U1 = runif(1e7)
2 > U2 = runif(1e7)
3 > R = sqrt(-2*log(U1))
4 > Theta = 2*pi*U2
5 > Z = R*cos(Theta)
6 > X = 2+3*Z
7 > hist(X, proba=TRUE)
8 > curve(dnorm(x ,2 ,3) ,add=TRUE)
```



Génération de loi Gaussienne $\mathcal{N}(\mu, \sigma^2)$ censurée

On veut simuler X conditionnellement à $X > 0$

```
1 > U1 = runif(1e7)
2 > U2 = runif(1e7)
3 > R = sqrt(-2*log(U1))
4 > Theta = 2*pi*U2
5 > Z = R*cos(Theta)
6 > X = 2+3*Z
7 > X = X[X>0]
8 > hist(X, proba=TRUE)
```

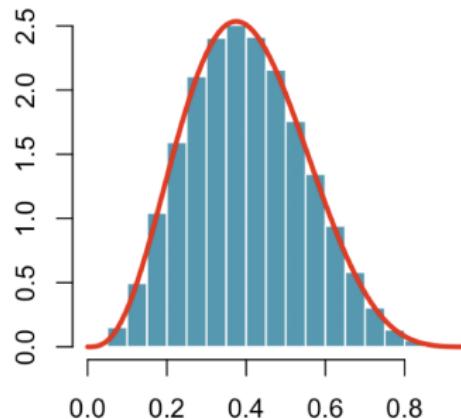


Génération de loi Gamma

C'est plus compliqué...

On peut utiliser les fonctions R pour les lois usuelles,
`runif`, `rbinom`, `rpois`, `rexp`, `rnorm`, `rlnorm`, `rgamma`, etc.

```
1 > a = 4
2 > b = 6
3 > Z = rgamma(1e7, a, b)
4 > hist(Z, probability=TRUE)
5 > curve(dgamma(x,a,b))
```



Notion de mélange

Sommme de fonctions de répartition

Si F_1 et F_2 sont deux fonctions de répartition, si $p_1, p_2 \in (0, 1)$, avec $p_1 + p_2 = 1$

$$F : x \mapsto p_1 F_1(x) + p_2 F_2(x)$$

est également une fonction de répartition

Mélange de lois

Si $X_1 \sim F_1$ et $X_2 \sim F_2$,

$$X = \mathbf{1}(Z = 1)X_1 + \mathbf{1}(Z = 2)X_2 = \begin{cases} X_1 & \text{si } Z = 1 \text{ (proba } p_1) \\ X_2 & \text{si } Z = 2 \text{ (proba } p_2) \end{cases}$$

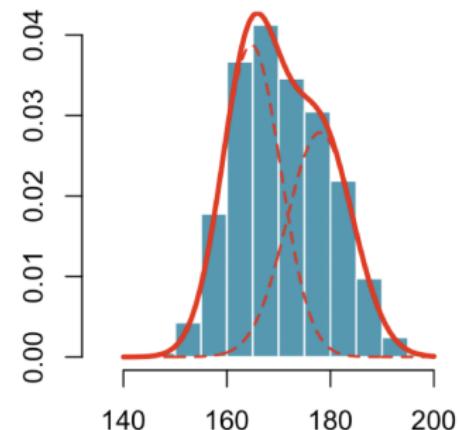
a pour loi $F(x) = p_1 F_1(x) + p_2 F_2(x)$.

Génération d'une loi mélange

Soit $I \sim \mathcal{B}(p)$, $p \in (0, 1)$

$$X \sim \begin{cases} \mathcal{N}(\mu_1, \sigma_1^2) & \text{si } I = 1 \\ \mathcal{N}(\mu_2, \sigma_2^2) & \text{si } I = 2 \end{cases}$$

```
1 > m1 = 178
2 > m2 = 164
3 > s1 = 6.44
4 > s2 = 5.66
5 > p = 0.45
6 > I = sample(1:2, size = 1e6,
               prob = c(p, 1-p), replace =
               TRUE)
7 > Z = rnorm(1e6, m1, s1)*(I==1) +
               rnorm(1e6, m2, s2)*(I==2)
8 > hist(Z, proba=TRUE)
```



Génération d'un vecteur Gaussien ★★★

Décomposition de Cholesky

Soit Σ une matrice de variance-covariance (matrice symétrique définie positive), il existe une matrice triangulaire inférieure L telle que $\Sigma = LL^\top$,

$$L = \begin{bmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & \vdots & \ddots & \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{bmatrix}$$

Il existe une unique matrice triangulaire inférieure dont les termes de la diagonale sont positifs telle que $\Sigma = LL^\top$.

Génération d'un vecteur Gaussien ★★★

Example:

$$\Sigma = \begin{bmatrix} 4 & -6 & 8 & 2 \\ -6 & 10 & -15 & -3 \\ 8 & -15 & 26 & -1 \\ 2 & -3 & -1 & 62 \end{bmatrix} \text{ et } L = \begin{bmatrix} 2 & 0 & 0 & 0 \\ -3 & 1 & 0 & 0 \\ 4 & -3 & 1 & 0 \\ 1 & 0 & -5 & 6 \end{bmatrix}$$

```
1 > S = matrix(c(4,-6,8,2,-6,10,-15,-3,8,-15,26,
   -1,2,-3,-1,62),4,4)
2 > chol(S)
3      [,1]  [,2]  [,3]  [,4]
4 [1,]     2    -3     4     1
5 [2,]     0     1    -3     0
6 [3,]     0     0     1    -5
7 [4,]     0     0     0     6
```

Génération d'un vecteur Gaussien ★★★

Vecteur Gaussien, $\mathcal{N}(\mu, \Sigma)$ et $\mathcal{N}(\mathbf{0}, \mathbb{I})$

$\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$ si et seulement si $\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}^T \mathbf{Z}$ où $\Sigma = \mathbf{L}\mathbf{L}^T$ et $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$.

$$\begin{cases} X_1 &= \mu_1 + L_{11}Z_1 \\ X_2 &= \mu_2 + L_{21}Z_1 + L_{22}Z_2 \\ X_3 &= \mu_3 + L_{31}Z_1 + L_{32}Z_2 + L_{33}Z_3 \\ \vdots \\ X_d &= \mu_d + L_{d1}Z_1 + L_{d2}Z_2 + \cdots + L_{d(d-1)}Z_{d-1} + L_{dd}Z_d \end{cases}$$