



STT 1000 - STATISTIQUES

ARTHUR CHARPENTIER



Random Numbers?

TABLE OF RANDOM DIGITS

1

00000	10097	32533	76520	13586	34673	54876	80959	09117	39292	74945
00001	37542	04805	64894	74296	24805	24037	20636	10402	00822	91665
00002	08422	68953	19645	09303	23209	02560	15953	34764	35080	33606
00003	99019	02529	09376	70715	38311	31165	88676	74397	04436	27659
00004	12807	99970	80157	36147	64032	36653	98951	16877	12171	76833
00005	66065	74717	34072	76850	36697	36170	65813	39885	11199	29170
00006	31060	10805	45571	82406	35303	42614	86799	07439	23403	09732
00007	85269	77602	02051	65692	68665	74818	73053	85247	18623	88579
00008	63573	32135	05325	47048	90553	57548	28468	28709	83491	25624
00009	73796	45753	03529	64778	35808	34282	60935	20344	35273	88435
00010	98520	17767	14905	68607	22109	40558	60970	93433	50500	73998
00011	11805	05431	39808	27732	50725	68248	29405	24201	52775	67851
00012	83452	99634	06288	98083	13746	70078	18475	40610	68711	77817
00013	88685	40200	86507	58401	36766	67951	90364	76493	29609	11062
00014	99594	67348	87517	64969	91826	08928	93785	61368	23478	34113
00015	65481	17674	17468	50950	58047	76974	73039	57186	40218	16544
00016	80124	35635	17727	08015	45318	22374	21115	78253	14385	53763
00017	74350	99817	77402	77214	43236	00210	45521	64237	96286	02655
00018	69916	26803	66252	29148	36936	87203	76621	13990	94400	56418
00019	09893	20505	14225	68514	46427	56788	96297	78822	54382	14598
00020	91499	14523	68479	27686	46162	83554	94750	89923	37089	20048
00021	80336	94598	26940	36858	70297	34135	53140	33340	42050	82341
00022	44104	81949	85157	47954	32979	26575	57600	40881	22222	06413
00023	12550	73742	11100	02040	12860	74697	96644	89439	28707	25815
00024	63606	49329	16505	34484	40219	52563	43651	77082	07207	31790
00025	61186	90446	26457	47774	51894	32799	55394	50602	40589	50527

Source [A Million Random Digits with 100,000 Normal Deviates](#),
RAND, 1955.

Random Numbers?

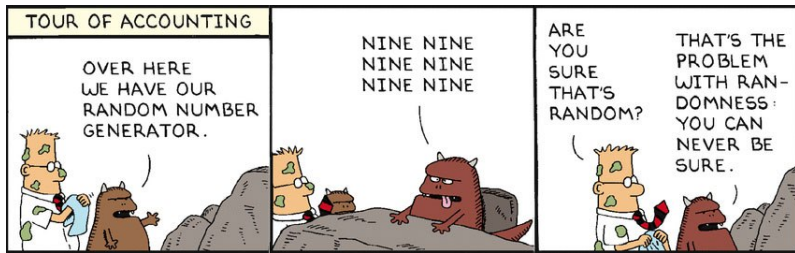
Here **random** means a sequence of numbers do not exhibit any discernible pattern, i.e. successively generated numbers can not be predicted.

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, quoted in **Knuth (1997)**

The goal of Pseudo-Random Numbers Generators is to produce a sequence of numbers in $[0, 1]$ that imitates ideal properties of random number.

```
1 > runif(50)
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
```

Random Numbers?



Source **Dibert, 2001.**

Randomness

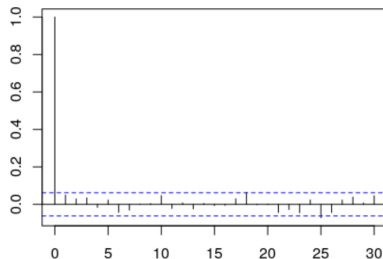
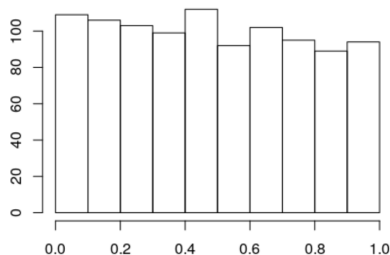
Heuristically,

1. calls should provide a **uniform sample**:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{u_i \in (a,b)} = b - a \text{ with } b > a,$$

2. calls should be **independent**: for $b > a$ and $d > c$.

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



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 $u_i = x_i/m$.

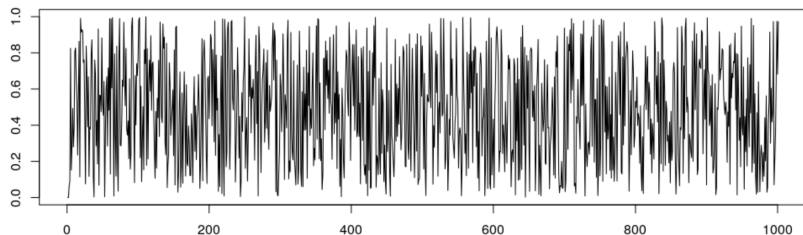
```
1 > a = 13; b = 43; m = 100; x = 77; u = rep(NA, 40)
2 > for (i in 1:40) {x = (a * x + b) %% m
3 +   u[i] = x / m }
4 > u
5 [1] 0.44 0.15 0.38 0.37 0.24 0.55 0.58 0.97 0.04 0.95
6 [11] 0.78 0.57 0.84 0.35 0.98 0.17 0.64 0.75 0.18 0.77
7 [21] 0.44 0.15 0.38 0.37 0.24 0.55 0.58 0.97 0.04 0.95
```

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

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

How to create randomness?

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



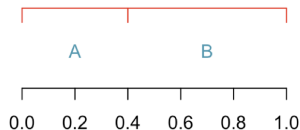
See [L'Ecuyer \(2017\)](#) for an historical perspective,

Note See [McCullough & Heiser \(2008\)](#) or [Mélard \(2014\)](#) about MS Excel and randomness

Génération de loi binomiale

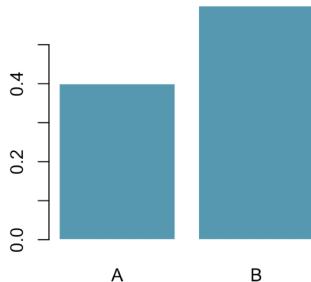
Soit $p \in (0, 1)$, $U \sim \mathcal{U}_{[0,1]}$

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



```
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,
  prob=c(1-p,p))
2
3 > table(Z)
4 Z
5      0      1
6 6000144 3999856
```



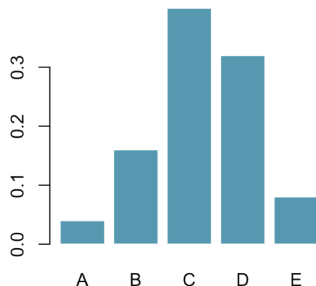
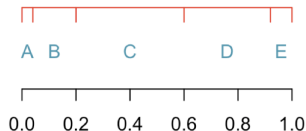
Génération de loi multinomiale

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



Inversion de la Fonction de Répartition

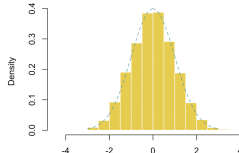
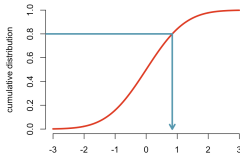
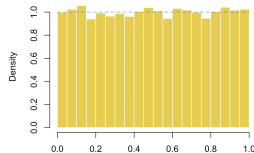
Inversion de fonction de répartition, **inverse method sampling**

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

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

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

where $F^{-1}(u) = \inf \{x \mid F(x) \geq u\}$ for $u \in (0, 1)$.



Inversion de la Fonction de Répartition

```
1 > U = runif(100)
2 [1] 0.26 0.35 0.31 0.76 0.52 0.06 0.03 0.23 0.67 0.14
3 [11] 0.17 0.13 0.58 0.93 0.32 0.11 0.53 0.13 0.09 0.19
4 [21] 0.32 0.37 0.91 0.47 0.28 0.38 0.88 0.98 0.49 0.84
5 [31] 0.51 0.63 0.14 0.60 0.79 0.17 0.37 0.33 0.46 0.72
6 [41] 0.92 0.39 0.42 0.48 0.70 0.30 0.05 0.51 0.38 0.27
7 [51] 0.51 0.69 0.21 0.11 0.17 0.19 0.14 0.68 0.99 0.50
8 [61] 0.26 0.69 0.43 0.25 0.06 0.26 0.32 0.10 0.18 0.08
9 [71] 0.05 0.55 0.13 0.50 0.75 0.18 0.15 0.12 0.81 0.35
```

```
1 > Q(U)
2 [1] 1.04 -0.48 0.81 -0.86 -0.33 0.74 0.92 0.38
3 [9] -0.80 0.95 -0.76 0.22 0.44 0.77 0.25 -1.45
4 [17] 0.10 -0.12 -1.87 0.68 0.73 -1.06 -0.19 -0.19
5 [25] -1.10 -0.48 1.09 1.11 0.06 0.04 0.15 0.08
6 [33] -0.45 -1.29 0.48 -0.33 0.95 0.25 0.80 1.58
7 [41] 0.31 -1.51 1.57 0.84 0.07 0.01 -0.96 0.56
8 [49] -0.66 0.49 0.46 -1.57 0.00 -0.29 1.89 0.60
9 [57] 0.34 0.43 1.01 0.31 -0.20 -0.19 -0.07 -0.07
10 [65] -0.04 1.31 -0.35 -0.37 -0.35 -2.26 1.47 -1.17
```

Inversion de la Fonction de Répartition Empirique

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

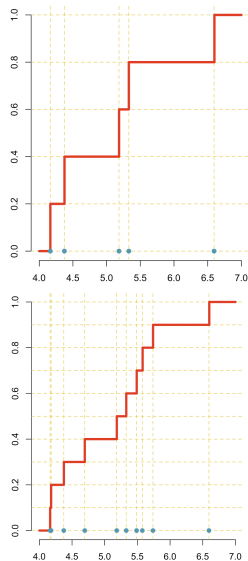
$$F(x) = \mathbb{P}[X \leq x],$$

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 $n \rightarrow \infty$,
or more precisely, almost surely

$$\|\hat{F}_n - F\|_{\infty} = \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \rightarrow 0$$



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
```

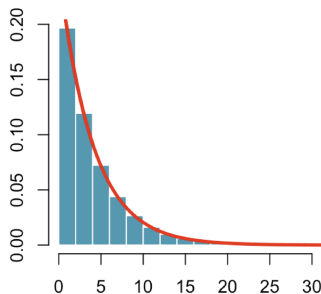
called **bootstrapping**

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 Exponentielle

Autre preuve? Soit U une loi uniforme, de densité $f(x) = \mathbf{1}_{[0,1]}(x)$, et considérons la fonction $g(x) = -b \log(x)$, telle que $h(y) = g^{-1}(y) = \exp[-y/b]$. Soit $Y = g(X)$, prenant les valeurs dans $(0, \infty)$.

Comme $h'(y) = -\exp[-y/b]/b$, comme $g(y) = f(h(y)) \cdot |h'(y)|$,

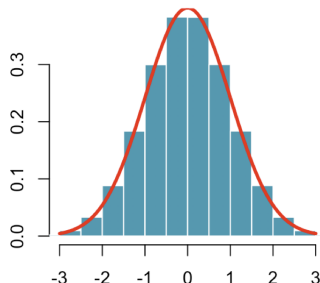
$$g(y) = f_X \left(e^{-y/b} \right) |h'(y)| = \frac{1}{b} e^{-y/b} \text{ sur } (0, \infty)$$

qui est la loi exponentielle de paramètre $a = 1/b$.

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)
```

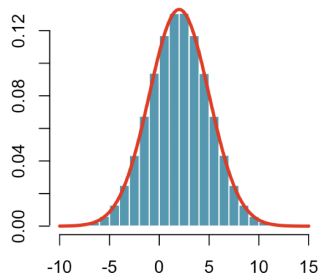


On parle de la méthode de Box-Muller.

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

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

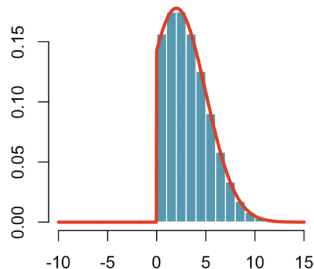
```
1 > U1 = runif(1e7)
2 > U2 = runif(1e7)
3 > R = sqrt(-2*log(U1))
4 > Theta = 2*pi*U2
5 > Z = R*cos(Theta)
6 > X = 2+3*Z
7 > hist(X, proba=TRUE)
8 > curve(dnorm(x,2,3),add=TRUE)
```



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)
```

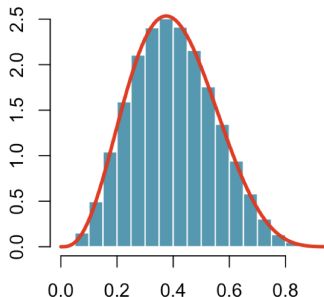


Génération de loi Gamma

C'est plus compliqué...

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

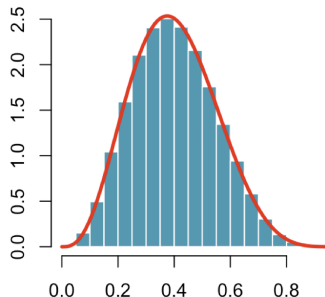
```
1 > a = 4
2 > b = 6
3 > Z = rgamma(1e7, a, b)
4 > hist(Z, probability=TRUE)
5 > curve(dgamma(x,a,b))
```



Génération de loi Beta

Si $U_1, U_2 \sim \mathcal{U}_{[0,1]}$, indépendantes, soient $V_1 = U_1^{1/a}$
 $V_2 = U_1^{1/a} + U_2^{1/b}$, alors V_1/V_2 sachant $V_1 \leq 1$ suit une loi $\mathcal{B}(a, b)$

```
1 > a = 4
2 > b = 6
3 > U1 = runif(1e7)
4 > U2 = runif(1e7)
5 > V1 = U1^(1/a)
6 > V2 = U1^(1/a)+U2^(1/b)
7 > Z = (V1/V2)[V2<=1]
8 > hist(Z, probability=TRUE)
9 > curve(dbeta(x,a,b),add=TRUE)
```

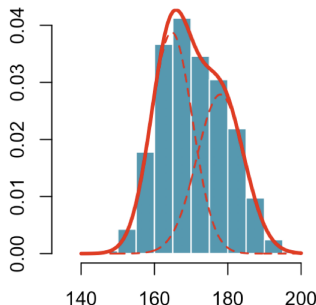


Génération d'une loi mélange

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

$$X \sim \begin{cases} \mathcal{N}(\mu_0, \sigma_0^2) & \text{si } I = 0 \\ \mathcal{N}(\mu_1, \sigma_1^2) & \text{si } I = 1 \end{cases}$$

```
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 $\Sigma = LL^T$,

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

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

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,  
2     -1,2,-3,-1,62),4,4)  
3  
4     [,1] [,2] [,3] [,4]  
5 [1,]    2   -3    4    1  
6 [2,]    0    1   -3    0  
7 [3,]    0    0    1   -5  
8 [4,]    0    0    0    6
```

Génération d'un vecteur Gaussien

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

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

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

Calcul d'espérance

$X \sim LN(0, 1)$, que vaut $\mathbb{P}[X > 3]$?

1. Calcul intégral

$$\mathbb{P}[X > 3] = \int_3^{\infty} \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{(\ln x - \mu)^2}{2\sigma^2}\right) dx$$

2. Calcul numérique de l'intégrale

```
1 > integrate(dlnorm,3,Inf)
2 0.1359686 with absolute error < 2.8e-05
```

3. Calcul numérique par simulations

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}(x_i > 3) \quad \text{où } x_i \leftarrow \text{loi log-normale } LN(0, 1)$$

```
1 > mean(rlnorm(1e6)>3)
2 [1] 0.1359608
```

Calcul d'espérance (ou pas)

Example Loi de Pareto, $F(x) = 1 - \left(\frac{x_m}{x}\right)^\alpha$ pour $x \geq x_m$

Un algorithme simple pour la simuler est par l'inverse de la fonction de répartition, $T = \frac{x_m}{U^{1/\alpha}}$ xxx

Visualisation de lois

Nous avons vu que si U_1, \dots, U_n est une collection de variables uniformes indépendantes, et si $U_{(k)}$ désigne la k ième observation, $U_k \sim \mathcal{B}(k, n - k + 1)$. Aussi

$$\min\{U_1, \dots, U_n\} \sim \mathcal{B}(1, n) \text{ et } \max\{U_1, \dots, U_n\} \sim \mathcal{B}(n, 1)$$

```
1 > n = 10
2 > ns = 1e4
3 > U = matrix(runif
      (n*ns), ns, n)
4 > minU = apply(U
      , 1, min)
5 > maxU = apply(U
      , 1, max)
6 > hist(minU)
```

