

OS13 : Test d'estimation et loi extrême des
lois de Cauchy et géométrie

ADRIEN WARTELLE
TRAN QUOC NHAT HAN

3 octobre 2018

Sommaire

1	Loi de Cauchy	1
1.1	Rappel	1
1.2	Estimer les paramètres inconnus	2
1.3	Test de paramètres	4
1.4	Test d'adéquation	5
1.5	Etude de biais d'estimateur	5
1.6	Etude de loi d'extremum	6
2	Loi géométrique	11
2.1	Génération de n variables aléatoires	12
2.2	Estimateurs du maximum de vraisemblance	13

Résumé

Ce rapport est à montrer l'application de certaines tests statistiques pour la loi géométrique et la loi de Cauchy.

1 Loi de Cauchy

1.1 Rappel

Soit f_X la fonction de densité de Cauchy de deux paramètres x_0 et a ($a > 0$), définie par :

$$f_X(x) = \frac{1}{\pi a \left(1 + \left(\frac{x-x_0}{a}\right)^2\right)} = \frac{1}{\pi} \frac{a}{(x-x_0)^2 + a} \quad (1.1)$$

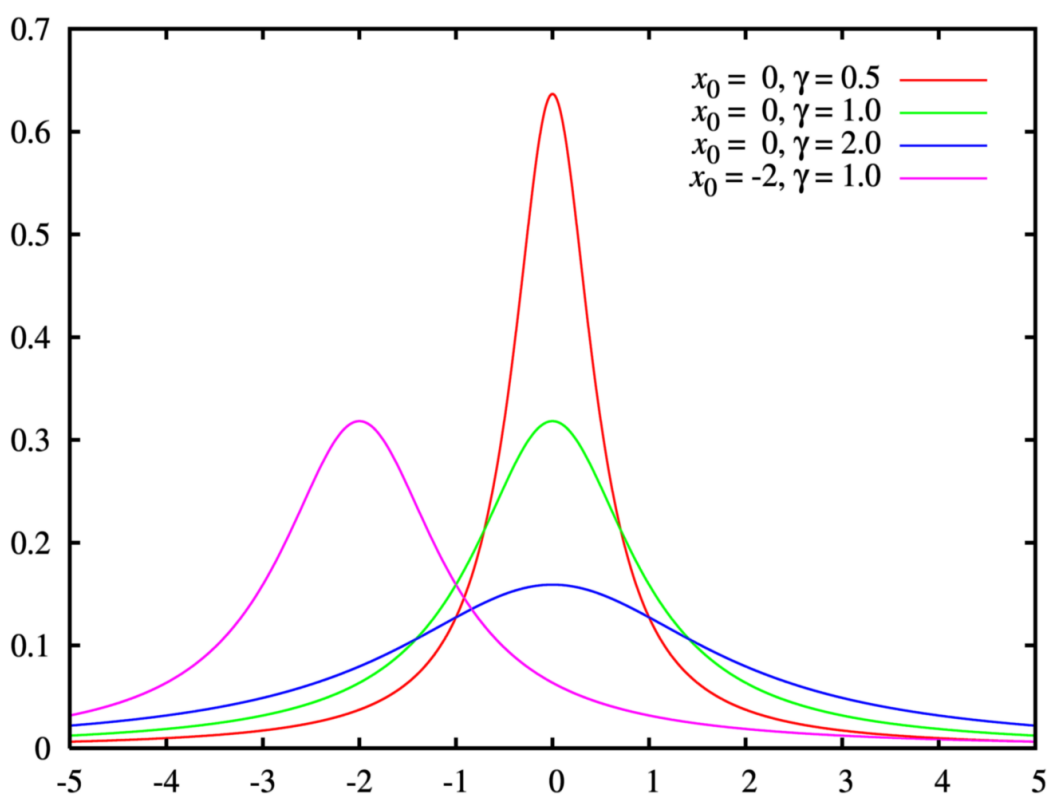


FIGURE 1 – Distribution théorique de la loi Cauchy. Source : [1]

a est dite l'échelle de la fonction, et x_0 est son médian.

La loi de Cauchy n'admet ni espérance ni écart type.

La fonction de répartition :

$$F_X(x) = \frac{1}{\pi} \arctan\left(\frac{x-x_0}{a}\right) + \frac{1}{2} \quad (1.2)$$

1.2 Estimer les paramètres inconnus

Calcul théorique

Soient n réalisations x_1, x_2, \dots, x_n . Assumons que ces données suivent la loi de Cauchy (1.1).

Nous allons utiliser la méthode de maximum de rapport de vraisemblance.

En l'absence d'information de la distribution, nous assumons que ces mesures sont indépendants. Posons une variable aléatoire X_i correspondante à chaque réalisation x_i pour $i = \overline{1, n}$.

On écrit la loi conjointe de ces n variables :

$$\begin{aligned} L(X) &= L(X_1, X_2, \dots, X_n) = \prod_{i=1}^n f_X(X_i) \\ &= \prod_{i=1}^n \frac{1}{\pi} \frac{a}{(x_i - x_0)^2 + a} \\ &= \pi^{-n} \prod_{i=1}^n \frac{a}{(x_i - x_0)^2 + a} \end{aligned}$$

On cherche l'optimum maximale.

$$\begin{aligned} \frac{\partial L}{\partial x_0}(X) &= \pi^{-n} \sum_{j=1}^n \left(-\frac{a(2x_j - 2x_0)}{[(x_0 - x_j)^2 + a]^2} \prod_{i=1; i \neq j}^n \frac{a}{(x_i - x_0)^2 + a} \right) \\ &= \pi^{-n} \left(\prod_{i=1}^n \frac{a}{(x_i - x_0)^2 + a} \right) \left(-\sum_{j=1}^n \frac{2x_0 - 2x_j}{(x_0 - x_j)^2 + a} \right) \\ \frac{\partial L}{\partial x_0}(X) &= 0 \Leftrightarrow \sum_{j=1}^n \frac{x_0 - x_j}{(x_0 - x_j)^2 + a} = 0 : \text{insolvable par la main} \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial a}(X) &= \pi^{-n} \sum_{j=1}^n \left(\frac{(x_j - x_0)^2 + a - a}{(x_j - x_0)^2 + a} \prod_{i=1; i \neq j}^n \frac{a}{(x_i - x_0)^2 + a} \right) \\ &= \pi^{-n} a^{n-1} \frac{\sum_{j=1}^n (x_0 - x_j)^2}{\prod_{i=1}^n ((x_i - x_0)^2 + a)} > 0 \forall a > 0 \end{aligned}$$

Résolution par R

On génère une échantillon avec a et x_0 au choix. Puis, utiliser la fonction *mle* (Maximum Likelihood Estimator) de librairie *stats4* avec 2 valeurs initiales \hat{a} et \hat{x}_0 pour estimer a et x_0 .

Nous essayons d'estimer x_0 et a directement. L'algorithme d'approximation implémenté dans R a besoin un bon point de départ, sinon le résultat obtenu variera grossièrement.

- Etant donné que la médiane est théoriquement aussi x_0 , choisissons x_0 comme la médiane de l'échantillon.
 - Nous avons $f_X(x_0) = \frac{1}{\pi a}$. En plus, $f_X(x_0)$ est la valeur maximale d de densité. Prenons alors $\hat{a} = \frac{1}{\pi d}$.
- Testons avec $x_0 = 13$ et $a = 0.5$.

```
1 # fixer le gemme et paramètres réels
2 set.seed(2018)
3 real_x0 = 13
4 real_a = 0.5
5 nbreaks = 40
6
7 # générer échantillons
8 N = 1000
9 x = rcauchy(n = N,
10            location = real_x0,
11            scale = real_a)
12
13 # définir fonction log négative pour MLE
14 ll = function(location, scale) {
15     dist = suppressWarnings(dcauchy(x,
16                                     location,
17                                     scale,
18                                     log = TRUE))
19     -sum(dist)
20 }
21
22 # MLE (Maximum Likelihood Estimator)
23 library(stats4)
24 hat_x0 = median(x)
25 d = hist(x,
26          breaks = nbreaks,
27          plot = FALSE) # histogram
28 hat_a = 1 / (pi * max(d[["density"]]))
29 result = mle(ll,
30             start = list(location = hat_x0, scale = hat_a))
31
32 # illustrer
33 hist(x,
34      freq = F,
35      breaks = nbreaks,
36      col = "green",
37      xlab = "x",
38      ylab = "Densité",
39      ylim = c(0, 0.2),
40      main = "Loi de Cauchy - réalité vs théorique") # histogram
41 curve(dcauchy(x, real_x0, real_a),
42       add = TRUE,
```

```

43     col = "black") # theoretic
44 curve(dcauchy(x, result@coef[1], result@coef[2]),
45       add = TRUE,
46       col = "red") # calculated
47 legend("topright",
48       legend = c("histogramme", "théorique", "réalité"),
49       fill = c("green", "black", "red"))

```

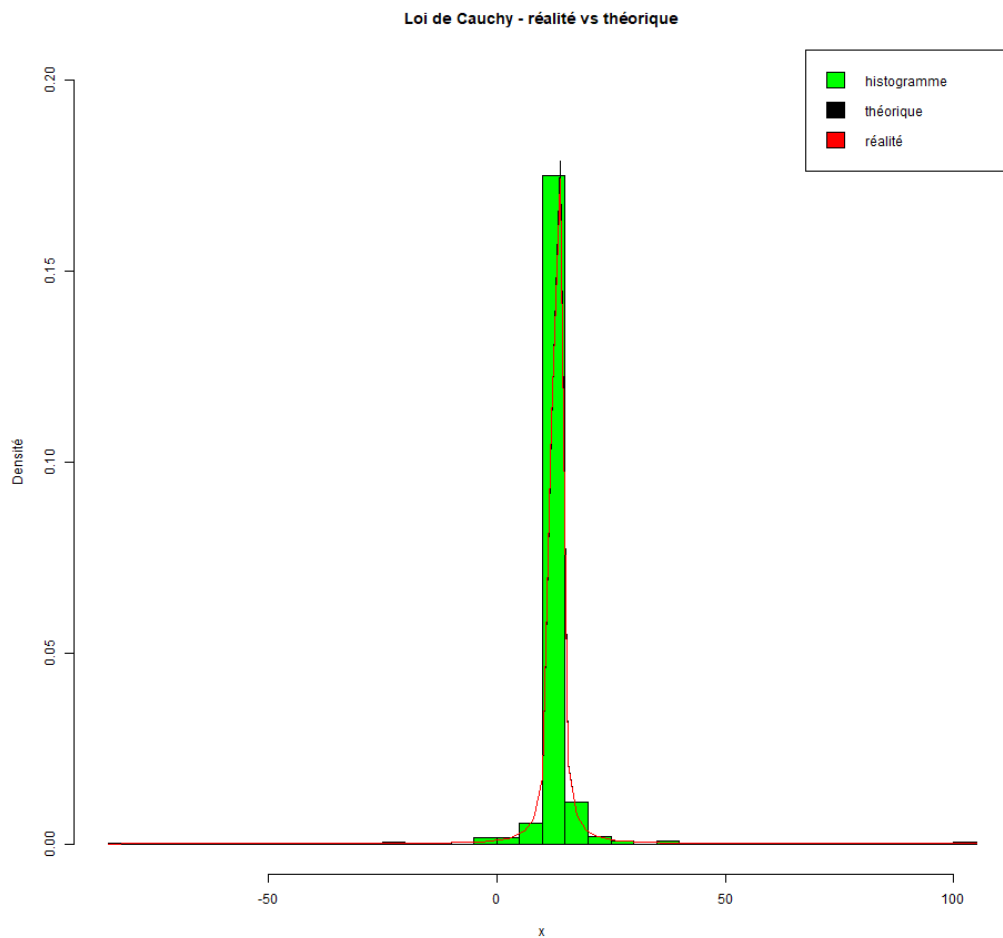


FIGURE 2 – Représentation de l’histogramme et de la courbe de densité

Pour $x_0 = 13$ et $a = 0.5$, les valeurs trouvées par R sont $\hat{x}_0 = 12.9877139$ et $\hat{a} = 0.4951073$.

1.3 Test de paramètres

Nous testons ici si les valeurs trouvées par R sont vraiment approchées à celles théoriques. Pour l’instant, nous n’avons pas d’outils pour vérifier 2

variables en couple.

1.4 Test d'adéquation

R nous donne la fonction *ks.test* pour valider la compabilité d'une échantillon avec une loi selon le test de Kolmogorov-Smirnov.

```
1 ks.test(x, "pcauchy", real_x0, real_a)
```

On a trouvé $D = 0.02696$, $p - value = 0.4614$, signifiant l'écart maximale est 0.02696 et le niveau d'acceptance est 0.4614. Vu que nous fixons $\alpha = 5\% < p - value$, l'échantillon dépasse largement le test. Autrement dit, il se distribue selon la loi de Cauchy.

1.5 Etude de biais d'estimateur

Maintenant on s'intéresse au biais d'estimateur. On relance l'algorithme en-dessus avec de différents observations ($N = 100, 125, 150, \dots, 1100$). $x_0 = 13$ et $a = 0.5$.

```
1 library(stats4)
2 # définir fonction log négative pour MLE
3 ll = function(location, scale) {
4     dist = suppressWarnings(dcauchy(x,
5                                     location,
6                                     scale,
7                                     log = TRUE))
8     - sum(dist)
9 }
10
11 # fixer les paramètres réels
12 set.seed(2018)
13 real_x0 = 13
14 real_a = 0.5
15 nbreaks = 40
16 minN = 100
17 numN = 40
18 delta = 25
19 maxN = minN + (numN - 1) * delta
20
21 Ns = seq.int(minN, maxN, delta) # tous 40 valeurs de n de 100 à 1000
22 results = array(NA, c(2, numN)) # à contenir tous les résultats
23
24 for (i in 1:numN) {
25     N = minN + (i - 1) * delta
26
27     # générer échantillon
28     x = rcauchy(n = N,
29                location = real_x0,
30                scale = real_a)
31
32     # MLE (Maximum Likelihood Estimator)
33     hat_x0 = median(x)
34     d = hist(x,
35              breaks = nbreaks,
```



```

36         plot = FALSE) # histogram
37     hat_a = 1 / (pi * max(d[["density"]]))
38     result = mle(ll,
39                 start = list(location = hat_x0, scale = hat_a))
40
41     # récupérer le résultat
42     results[1, i] = result@coef[1] # x_0
43     results[2, i] = result@coef[2] # a
44 }
45
46 # illustrer
47 plot(x = results[1, ],
48      y = results[2, ],
49      col = rgb(0, 100, 0, 70, maxColorValue = 255),
50      pch = 16,
51      cex = 1.5,
52      main = "Convergence d'estimateurs",
53      xlab = "Location x_0",
54      ylab = "Echelle a"
55     )
56 points(x = real_x0,
57        y = real_a,
58        col = "red",
59        pch = 16,
60        cex = 2
61       )
62
63 legend("topleft",
64       legend = c("valeur théorique", "valeur estimée"),
65       fill = c("red", rgb(0, 100, 0, 70, maxColorValue = 255)))

```

Nous observons que les valeurs estimées se rassemblent assez proches de la valeur théorique.

1.6 Etude de loi d'extremum

Générons donc des échantillons de n valeurs ($n = 10, 100, 1000, \text{etc.}$). Parmi chacun, $\frac{n}{10}$ valeurs les plus grandes seront gardées pour étudier la loi d'extremum.

Comme les réalisations sont indépendantes et identiquement distribuées, théoriquement :

$$\begin{aligned}
 F_{x_{\max}}(x) &= P(X_{\max} \leq x) \\
 &= P(X_1 \leq x, \dots, X_n \leq x) \\
 &= \prod_{i=1}^n F_X(x) \\
 &= F_X(x)^n
 \end{aligned}$$

Avec $F_X(x) = \frac{1}{\pi} \arctan \frac{x-x_0}{a} + \frac{1}{2}$, la fonction de répartition de la loi de Cauchy. Alors,

$$F_{x_{\max}}(x) = \left(\frac{1}{\pi} \arctan \frac{x-x_0}{a} + \frac{1}{2} \right)^n$$

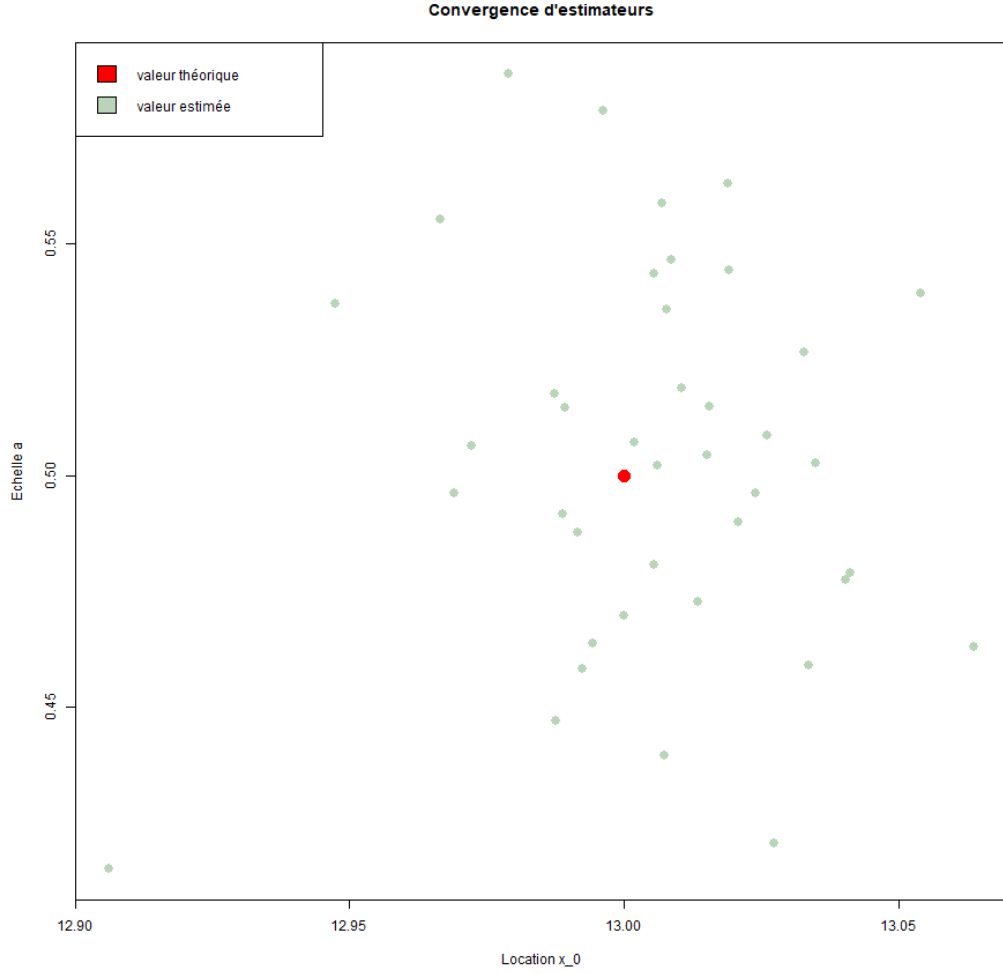


FIGURE 3 – Convergences de \hat{x}_0 et \hat{a} vers valeurs théoriques

Or,

$$\left| \arctan \frac{x - x_0}{a} \right| < \frac{\pi}{2} \Rightarrow \frac{-\pi}{2} < \arctan \frac{x - x_0}{a} < \frac{\pi}{2} \Rightarrow 0 < \arctan \frac{x - x_0}{a} + \frac{\pi}{2} < \pi$$

Par conséquence,

$$\begin{aligned} \tan \left(\arctan \frac{x - x_0}{a} + \frac{\pi}{2} \right) &= -\cot \left(\arctan \frac{x - x_0}{a} \right) = -\frac{a}{x - x_0} \\ \Rightarrow \arctan \frac{x - x_0}{a} + \frac{\pi}{2} &= \arctan \left(-\frac{a}{x - x_0} \right) = \arctan \frac{a}{x_0 - x} \\ \Rightarrow F_{x_{\max}}(x) &= \pi^{-n} \left(\arctan \frac{x - x_0}{a} + \frac{\pi}{2} \right)^n = \pi^{-n} \arctan^n \frac{a}{x_0 - x} \end{aligned}$$

Nous prenons 3 valeurs les plus grandes pour chaque itération, et établissons leur histogramme.

```

1  # fixer les paramètres réels
2  set.seed(2018)
3  real_x0 = 13
4  real_a = 0.5
5  minN = 1000
6  numN = 100
7  delta = 500
8  maxN = minN + (numN - 1) * delta
9  pickupNum = 3
10 results = rep(NA, 3000) # à contenir tous les résultats
11 countResults = 0
12
13 for (i in 1:numN) {
14   N = minN + (i - 1) * delta
15
16   # générer échantillon
17   x = rcauchy(n = N,
18             location = real_x0,
19             scale = real_a)
20
21   # trier
22   sort(x, decreasing = TRUE) # dans l'ordre décroissant
23
24   # tirer
25   extremes = head(x, pickupNum)
26   for (j in 1:length(extremes))
27     results[countResults + j] = extremes[j]
28   countResults = countResults + length(extremes)
29 }
30
31 # histogramme
32 hist(results,
33       freq = F,
34       breaks = 100,
35       col = "green",
36       xlab = "Extremum",
37       ylab = "Densité",
38       main = "Loi d'extremum au cas de Cauchy"
39 )

```

En fin d'approximer la loi d'extremum, nous ne regardons que les réalisations positives. En utilisant le package *poweRlaw*, nous calculons la compatibilité de la loi d'extremum.

```

1  results = na.omit(results) # filtrer les NA
2  library(poweRlaw)
3  positiveResults = results[which(results > 0)] # filtrer les négatifs
4  cauchy_pl = compl$new(positiveResults) # objet de calculer la loi d'extremum
5  est = estimate_xmin(cauchy_pl) # estimer la borne inférieure
6  cauchy_pl$setXmin(est) # mettre à jour l'objet de distribution
7  plot(cauchy_pl, col = "green")
8  lines(cauchy_pl, col = "red")
9  bs_p = bootstrap_p(cauchy_pl) # tester l'estimation
10 bs_p$p

```

Calcul de p – *value* nous retourne 0, qui signifie l'incompatibilité de ces deux lois.

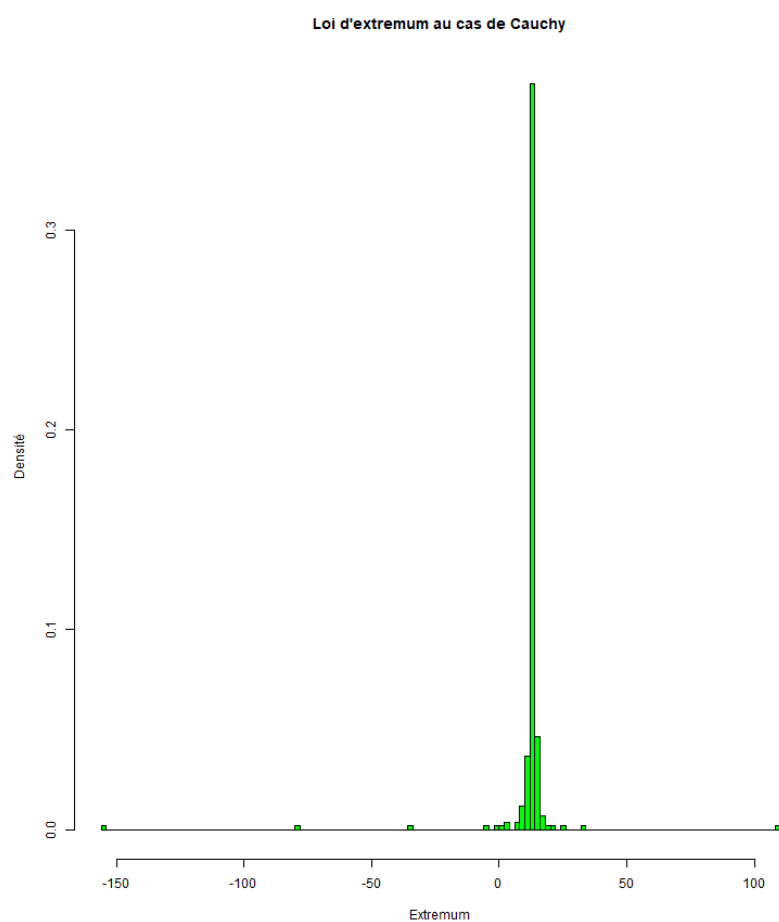


FIGURE 4 – Histogramme de valeurs extremum de la distribution de Cauchy

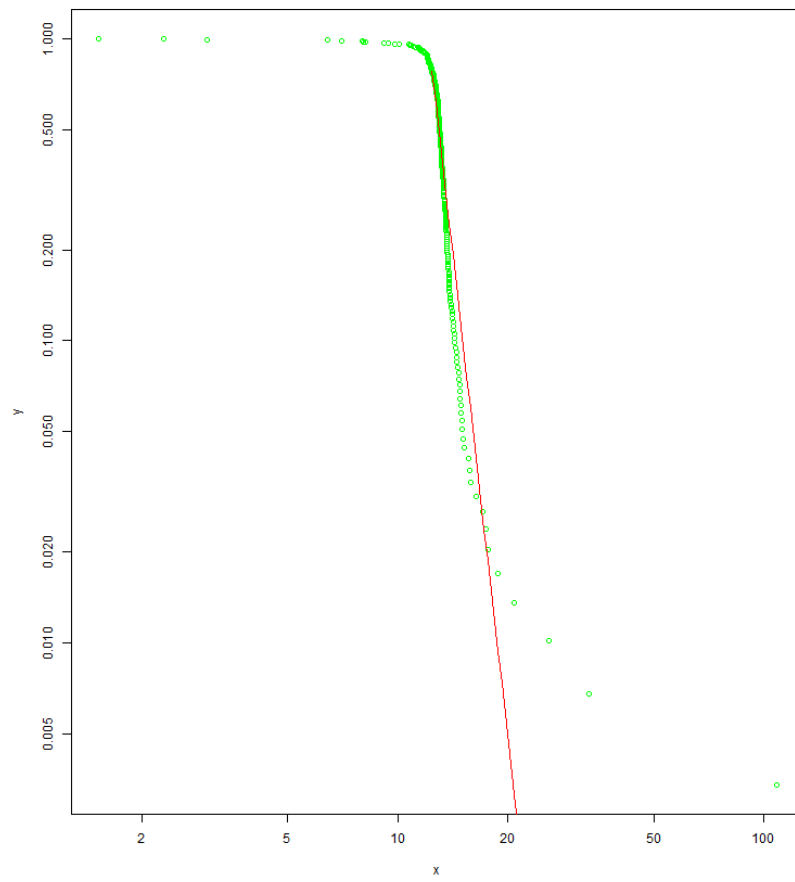


FIGURE 5 – Approcher la loi d'extremum avec celle des maximales de Cauchy. Les données (verte) ne s'approchent pas totalement l'approximation (rouge).

2 Loi géométrique

La loi géométrique est une loi discrète de distribution utilisé dans le cadre d'un enchainement (non fini) d'épreuves de Bernoulli. Soit X une variable aléatoire suivant cette loi, on a :

$$P(X = k) = (1 - p)^{k-1} p \quad k \in \mathbb{N}_+$$

La variable X est le numéro de la première épreuve où l'on obtient un succès qui a une probabilité p . Ainsi, p est le seul paramètre de loi (à estimer). La fonction de répartition $F(k) = P(X \leq k) \quad k \in \mathbb{N}_+$ est :

$$F(k) = 1 - (1 - p)^k$$

Démonstration. On peut remarquer que :

$$P(X = k) = (1 - p)^{k-1} (1 - (1 - p)) = (1 - p)^{k-1} - (1 - p)^k$$

Donc :

$$\begin{aligned} F(k) &= \sum_{l=1}^k P(X = l) = \sum_{l=1}^k ((1 - p)^{l-1} - (1 - p)^l) \\ F(k) &= \sum_{l=0}^{k-1} (1 - p)^l - \sum_{l=1}^k (1 - p)^l \end{aligned}$$

Soit $q = 1 - p$, on a :

$$\begin{aligned} F(k) &= \sum_{l=0}^{k-1} q^l - \sum_{l=1}^k q^l \\ F(k) &= \frac{1 - p^k}{1 - q} - \left(\frac{1 - q^{k+1}}{1 - q} - 1 \right) \\ F(k) &= \frac{1 - q^k - 1 + q^{k+1} + 1 - q}{1 - q} \\ F(k) &= \frac{1 - q + q^{k+1} - q^k}{1 - q} \\ F(k) &= \frac{1 - q - q^k(1 - q)}{1 - q} \\ F(k) &= \frac{1 - q - q^k(1 - q)}{1 - q} \\ F(k) &= 1 - q^k \end{aligned}$$

On retrouve bien :

$$F(k) = 1 - (1 - p)^k$$

□

2.1 Génération de n variables aléatoires

Afin de voir à quoi ressemble la distribution et de pouvoir tester l'estimateur que nous allons calculer par la suite, on génère des échantillons de 50, 500 et 5000 variables. On utilise ainsi le code Matlab (Octave) ci dessous :

```
1 p=0.1; %probability to win
2 pop1 = geomGen(50,p);
3 pop2 = geomGen(500,p);
4 pop3 = geomGen(5000,p);
5 k=1:max(pop3)+1;
6 distrib = (1-p).^k.*p;
7 subplot(221);
8 hist(pop1,k,1);
9 xlabel('k');
10 ylabel('Freq(k)');
11 title('Histogramme de 50 variables');
12 subplot(222);
13 hist(pop2,k,1);
14 xlabel('k');
15 ylabel('Freq(k)');
16 title('Histogramme de 500 variables');
17 subplot(223);
18 hist(pop3,k,1);
19 xlabel('k');
20 ylabel('Freq(k)');
21 title('Histogramme de 5000 variables');
22 subplot(224);
23 bar(k,distrib);
24 xlabel('k');
25 ylabel('P(X=k)');
26 title('Distribution théorique géométrique');
```

On obtient ainsi la figure 6.

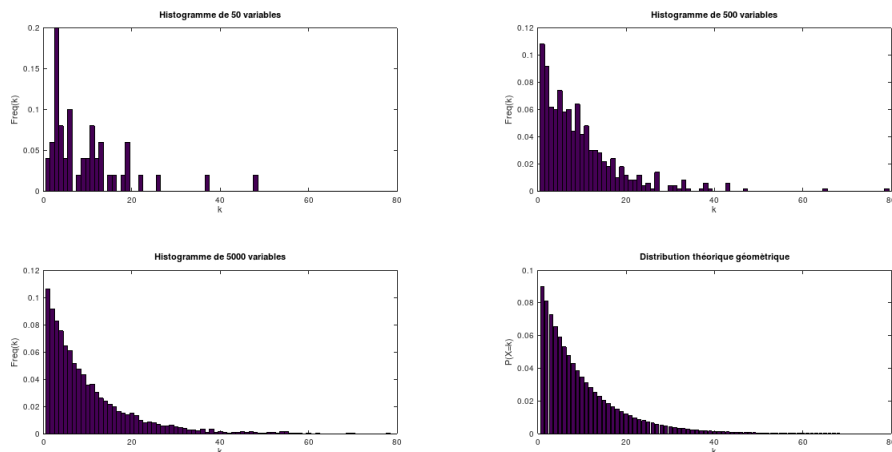


FIGURE 6 – Histogrammes et distribution de loi géométrique

Pour générer les populations (voir figure 6 , on utilise la fonction "geom-Gen" qui simule n expériences où, pour chacune d'entre elle, on effectue un

enchainement d'épreuve de Bernoulli en générant une variable de loi uniforme sur $[0; 1]$. Si la valeur obtenue est supérieure à p , on arrête la boucle et la valeur générée est égale au nombre d'épreuves, sinon on continue jusqu'à obtenir cette condition. Le code Matlab (Octave) utilisé est :

```

1 function [X]=geomGen(n,p)
2 X=ones(1,n);
3 for (i=1:n)
4     trial=rand(1,1);
5     k=1;
6     while (trial > p)
7         trial=rand(1,1);
8         k++;
9     endwhile
10    X(i)=k;
11 endfor
12 endfunction

```

2.2 Estimateurs du maximum de vraisemblance

L'estimateur \hat{p} du maximum de vraisemblance est la valeur qui maximise la loi de vraisemblance $L(p)$: $\hat{p} = \operatorname{argmax}(L(p))$. On calcule $L(p)$ dans un premier temps :

$$L(p) = L(X_1(p), X_2(p), \dots, X_n(p)) = \prod_{i=1}^n (1-p)^{X_i-1} p$$

avec X_i , $i \in \{1, \dots, n\}$ les variables d'échantillon.

$$L(p) = p^n (1-p)^{(\sum_{i=1}^n X_i) - n}$$

On peut effectuer un passage au logarithme pour trouver un maximum car il s'agit d'une fonction strictement croissante (et défini sur $]0; 1]$). L'argument du maximum du logarithme de $L(p)$ et du maximum de $L(p)$ sont les mêmes.

$$\log L(p) = n \log p + \left(\sum_{i=1}^n X_i \right) - n$$

Références

- [1] Wikipedia,
https://en.wikipedia.org/wiki/Cauchy_distribution