

# TP Statistiques : Compte-rendu

Michel Yoeung, Charles-Frédéric Amaudruz, Alexandre Berrada (ENSIMAG - 1A - Groupe 2)

Avril 2018

## 1 Première stratégie.

### Question 1

$$\forall i = 1, \dots, n, X_i = \begin{cases} 1 & \text{si le } i^{\text{eme}} \text{ poisson est bagué} \\ 0 & \text{sinon} \end{cases}$$

De plus, les résultats des pêches successives sont indépendants et la probabilité que le  $i^{\text{me}}$  poisson soit bagué est  $p = \frac{n_0}{\theta}$ .

Donc les  $X_i$  sont indépendants et identiques (iid) car suivent tous une loi de Bernoulli de paramètre de succès  $p = \frac{n_0}{\theta}$ .

Soit  $\theta = 1000$  et  $n_0 = 50$ . On choisit  $n = 60$  pour notre simulation sur R.

On obtient l'échantillon de données suivant :

```
[1] "échantillon :"  
[1] 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
[60] 0
```

On calcule ensuite sur R les moyennes et les variances (empiriques et théoriques) :

```
> source('~/.ensimag1/Statistiques/projet/premiere_strategie.R')  
[1] "moyenne empirique :"  
[1] 0.06666667  
[1] "moyenne théorique :"  
[1] 0.05  
[1] "variance empirique :"  
[1] 0.06222222  
[1] "variance empirique :"  
[1] 0.0475
```

On constate que pour ce jeu de données la moyenne empirique est assez proche de la moyenne théorique et de même pour les variances.

### Question 2

$T$  est une variable aléatoire qui compte le nombre de poissons bagués parmi les  $n$  poissons pêchés donc  $T = \sum_{i=1}^n X_i$  avec  $X_i$  qui suit une loi de Bernoulli.

Donc  $T$  suit une loi binomiale de paramètres  $n$  et  $p = \frac{n_0}{\theta}$ .

On donne ainsi avec R  $t$  le nombre de poissons pêchés sur notre échantillon de données :

```
[1] "nombre de poissons bagués parmi les n poissons pêchés :"  
[1] 4
```

### Question 3

*Estimateur des moments :*

Soit  $\overline{X_n}$  la moyenne empirique,

$$\begin{aligned} E[X] &\simeq \overline{X_n} \Rightarrow p \simeq \frac{1}{n} \sum_{i=1}^n x_i \\ &\Rightarrow \frac{n_0}{\theta} \simeq \frac{1}{n} \sum_{i=1}^n x_i \\ &\Rightarrow \theta \simeq \frac{n_0 * n}{\sum_{i=1}^n x_i} \end{aligned} \quad (1)$$

Donc l'estimateur des moments (d'ordre 1) vaut  $\tilde{\theta}_n = \frac{n_0 * n}{\sum_{i=1}^n x_i}$ .

*Estimateur de maximum de vraisemblance :*

$$\begin{aligned} L(X_1 = x_1, \dots, X_n = x_n, \theta) &= \prod_{i=1}^n P(X_i = x_i, \theta) \\ &= p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i} \\ &= \left(\frac{n_0}{\theta}\right)^{\sum_{i=1}^n x_i} \left(1 - \frac{n_0}{\theta}\right)^{n-\sum_{i=1}^n x_i} \end{aligned} \quad (2)$$

$$\ln(\mathcal{L}(X_1 = x_1, \dots, X_n = x_n, \theta)) = (\sum_{i=1}^n x_i)(\ln(n_0) - \ln(\theta)) + (n - \sum_{i=1}^n x_i) \ln(1 - \frac{n_0}{\theta})$$

En dérivant par rapport à  $\theta$  :

$$\begin{aligned} \frac{\partial \ln(\mathcal{L}(X_1 = x_1, \dots, X_n = x_n, \theta))}{\partial \theta} &= -\left(\sum_{i=1}^n x_i\right) \frac{1}{\theta} + \left(n - \sum_{i=1}^n x_i\right) \frac{\frac{n_0}{\theta^2}}{1 - \frac{n_0}{\theta}} \\ &= -\left(\sum_{i=1}^n x_i\right) \frac{1}{\theta} + \left(n - \sum_{i=1}^n x_i\right) \frac{n_0}{\theta^2 - \theta n_0} \\ &= 0 \Rightarrow \theta = \frac{n_0 * n}{\sum_{i=1}^n x_i} \end{aligned} \quad (3)$$

Donc l'estimateur de maximum de vraisemblance vaut  $\hat{\theta}_n = \frac{n_0 * n}{\sum_{i=1}^n x_i}$ .

Donc ces deux estimateurs sont confondus.

On calcule ainsi sur R la valeur de ces estimateurs sur notre échantillon de données :

```
##
[1] "estimateur des moments : "
[1] 750
[1] "estimateur de maximum de vraisemblance : "
[1] 750
#
```

Sur l'échantillon simulée, on observe un estimateur de  $\theta$  qui vaut 750 soit 25. Cet estimateur n'est pas très précis.

#### Question 4

intervalle de confiance exact de seuil  $\alpha$  pour  $\theta$  :

On a  $p = \frac{n_0}{\theta} \Rightarrow \theta = \frac{n_0}{p}$  donc en estimant  $p$  (avec la formule du cours), on obtient directement une estimation de  $\theta$ .

$$\left[ n_0 \left( 1 + \frac{n-T}{T+1} f_{2(n-T), 2(T+1), 1-\frac{\alpha}{2}} \right), n_0 \left( 1 + \frac{n-T+1}{T} f_{2(n-T+1), 2T, \frac{\alpha}{2}} \right) \right] \quad (4)$$

avec  $T = n\bar{X}_n = \sum_{i=1}^n x_i$  et  $f_{\nu_1, \nu_2, \alpha} = F_{\mathcal{F}}^{-1}(1 - \alpha, \nu_1, \nu_2)$ ,  $F_{\mathcal{F}}^{-1}$  étant ici la fonction quantile de la loi de Fisher-Snedecor.

intervalle de confiance asymptotique de seuil  $\alpha$  pour  $\theta$  :

On obtient l'intervalle associée avec la même démarche que précédemment.

$$\left[ \frac{n_0}{\bar{X}_n + u_\alpha \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}}, \frac{n_0}{\bar{X}_n - u_\alpha \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}} \right] \quad (5)$$

avec  $u_\alpha = F_{\mathcal{N}}^{-1}(1 - \frac{\alpha}{2})$ ,  $F_{\mathcal{N}}^{-1}$  étant ici la fonction quantile de la loi normale.

On calcule ainsi ces intervalles de confiance (exacts puis asymptotiques) associés aux valeurs de nos paramètres fixés sur  $R$  pour les différentes valeurs de  $\alpha$  spécifiées dans l'énoncé :

```
[1] "intervalle de confiance exact pour theta de seuil"
[1] 0.01
[1] 255.9594
[1] 4375.487
[1] "intervalle de confiance exact pour theta de seuil"
[1] 0.05
[1] 308.6672
[1] 2708.298
[1] "intervalle de confiance exact pour theta de seuil"
[1] 0.1
[1] 342.2379
[1] 2165.389
[1] "intervalle de confiance exact pour theta de seuil"
[1] 0.2
[1] 388.1135
[1] 1701.15
[1] .. .. .
```

On remarque que les intervalles exacts fournissent un encadrement plus précis de la vraie valeur de  $\theta$  que les intervalles asymptotiques ce qui semble logique car les intervalles asymptotiques sont plus efficaces pour un  $n$  très grand, or ici on a fixé  $n = 60$  seulement.

#### Question 5

$$P(\hat{\theta}_n = +\infty) = P\left(\frac{n_0 * n}{\sum_{i=1}^n x_i} = +\infty\right) = P\left(\sum_{i=1}^n x_i = 0\right) = P(T = 0) = (1 - p)^n = \left(1 - \frac{n_0}{\theta}\right)^n$$

Cet estimateur n'est pas convergent.

$$\text{Biais}(\hat{\theta}_n) = E[\hat{\theta}_n] - \theta$$

Comme  $\hat{\theta}_n$  a une probabilité non nulle de valoir  $+\infty$  alors  $E[\hat{\theta}_n] = +\infty \Rightarrow E[\hat{\theta}_n] - \theta = +\infty$  (car  $\theta$  est une constante)  $\Rightarrow \text{Biais}() = +\infty$ .

Donc on peut en déduire que le biais de cet estimateur vaut  $+\infty$ . Sur notre échantillon de données, on peut calculer la probabilité  $P(\hat{\theta}_n = +\infty)$  :

```

[1] "intervalle de confiance asymptotique pour theta de seuil"
[1] 0.01
[1] 334.1883
[1] -3070.703
[1] "intervalle de confiance asymptotique pour theta de seuil"
[1] 0.05
[1] 385.257
[1] 14085.18
[1] "intervalle de confiance asymptotique pour theta de seuil"
[1] 0.1
[1] 417.9345
[1] 3650.342
[1] "intervalle de confiance asymptotique pour theta de seuil"
[1] 0.2
[1] 463.2351
[1] 1968.753

[1] "probabilité que l'estimateur vale +infini :"
[1] 0.0460698

```

### Question 6

$$\begin{aligned}
 P(\hat{\theta}_n = +\infty) &> \frac{1}{2} \Rightarrow \left(1 - \frac{n_0}{\theta}\right)^n > \frac{1}{2} \\
 &\Rightarrow n \ln\left(1 - \frac{n_0}{\theta}\right) > \ln\left(\frac{1}{2}\right) \\
 &\Rightarrow n < -\frac{\ln(2)}{\ln\left(1 - \frac{n_0}{\theta}\right)} \quad (6) \\
 &\Rightarrow n \leq \left\lfloor -\frac{\ln(2)}{\ln\left(1 - \frac{n_0}{\theta}\right)} \right\rfloor
 \end{aligned}$$

Toujours avec notre échantillon de données, on calcule sur  $R$  cette valeur de  $n$  pour laquelle la probabilité  $P(\hat{\theta}_n = +\infty)$  :

```

[1] -----
[1] "pour que la probabilité que l'estimateur vale +infini soit strictement supérieure à 1/2, n doit être inférieur à :"
[1] 13
[1] "probabilité que l'estimateur vale +infini avec cette valeur de n :"
[1] 0.5133421

```

## 2 Deuxième stratégie.

Question 1

Question 2

Question 3

Question 4

Question 5

### 3 Application et comparaison des stratégies.

Question 1

Question 2

Question 3

Question 4

## 4 Vérifications expérimentales à base de simulations.

### Question 1

*On calcule  $R$  les différentes proportions d'appartenance de  $\theta$  pour les intervalles de confiance établis en variant les paramètres.*

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de theta (intervalles de confiance exacts) :"  
[1] 98 95 96 94 98 96  
[1] "proportions (en pourcentage) en fonction de l'augmentation de theta (intervalles de confiance asymptotiques) :"  
[1] 93 89 91 89 80 74
```

*Lorsqu'on augmente  $\theta$ , on remarque que la proportion d'appartenance aux intervalles de  $\theta$  diminue lorsqu'on choisit l'intervalle de confiance asymptotique.*

*En effet, lorsqu'on augmente  $\theta$ , la probabilité  $p$  diminue donc l'intervalle asymptotique a une plus grande probabilité de se "tromper" car le nombre  $n$  d'essais n'est pas grand.*

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de n0 (intervalles de confiance exacts) :"  
[1] 95 96 99 97 99 97  
[1] "proportions (en pourcentage) en fonction de l'augmentation de n0 (intervalles de confiance asymptotiques) :"  
[1] 87 94 88 97 95 89
```

*Lorsqu'on augmente  $n_0$ , on remarque que la proportion d'appartenance augmente très légèrement.*

*En effet, comme  $n_0$  est proportionnel à  $p$  et pour la même raison que précédemment, cela paraît cohérent.*

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de n (intervalles de confiance exacts) :"  
[1] 94 98 97 91 97 94  
[1] "proportions (en pourcentage) en fonction de l'augmentation de n (intervalles de confiance asymptotiques) :"  
[1] 95 96 97 96 96 93
```

*Lorsqu'on augmente  $n$ , on ne constate pas de changement particulier, ce qui peut paraître incohérent car concernant l'intervalle asymptotique, ce dernier devrait être plus précis lorsque  $n$  est grand.*

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de m (intervalles de confiance exacts) :"
```

[1]	95.00	94.00	95.20	94.60	95.02	95.62
-----	-------	-------	-------	-------	-------	-------

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de m (intervalles de confiance asymptotiques) :"
```

[1]	95.00	93.50	93.80	93.60	92.74	93.31
-----	-------	-------	-------	-------	-------	-------

Lorsqu'on augmente  $m$ , on constate une stabilisation au niveau des proportions donc de la précision des intervalles.

En effet, augmenter  $m$  permet juste de rendre la simulation plus précise puisqu'on se base sur un plus grand nombre d'essais pour établir les proportions.

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de alpha (intervalles de confiance exacts) :"
```

[1]	98	96	90	85	48	36
-----	----	----	----	----	----	----

```
[1] "proportions (en pourcentage) en fonction de l'augmentation de alpha (intervalles de confiance asymptotiques) :"
```

[1]	85	93	88	79	50	8
-----	----	----	----	----	----	---

Lorsqu'on augmente  $\alpha$ , quel que soit le type d'intervalle (exact ou asymptotique), il est cohérent que la proportion diminue car on baisse le niveau de confiance donc la précision de ces intervalles.

## Question 2

On simule  $m = 100$  échantillons de taille  $n = 5, 10, 100, 1000, 10000, 100000$  de loi de Bernoulli, puis on compare les  $m$  moyennes empiriques avec les espérances (avec une erreur de  $\epsilon = 0.01$ ) pour illustrer la loi faible des grands nombres :

```
[1] "valeur de n :"
```

[1]	5e+00	1e+01	1e+02	1e+03	1e+04	1e+05
-----	-------	-------	-------	-------	-------	-------

```
[1] "proportions (en pourcentage) en fonction de n :"
```

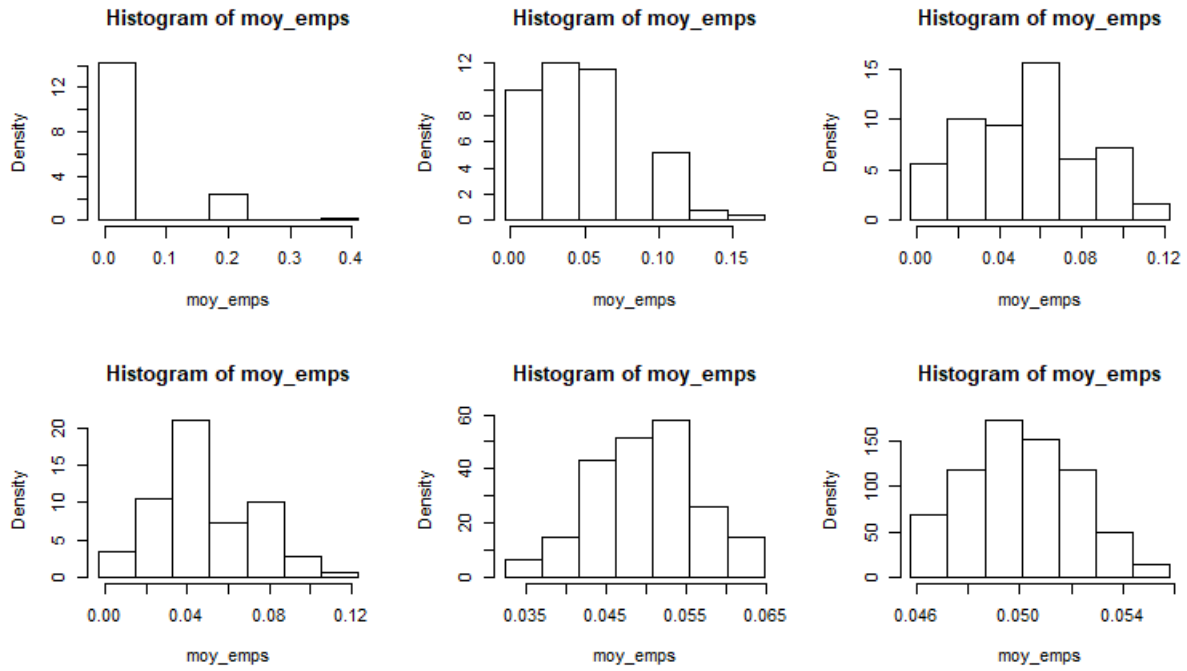
[1]	0	0	24	85	100	100
-----	---	---	----	----	-----	-----

Lorsque  $n$  augmente, on remarque que l'écart entre la moyenne empirique et l'espérance diminue. A partir de  $n = 1000$  environ, on peut dire que la moyenne empirique peut être approximée par l'espérance. La loi faible des grands nombres est ainsi illustrée.

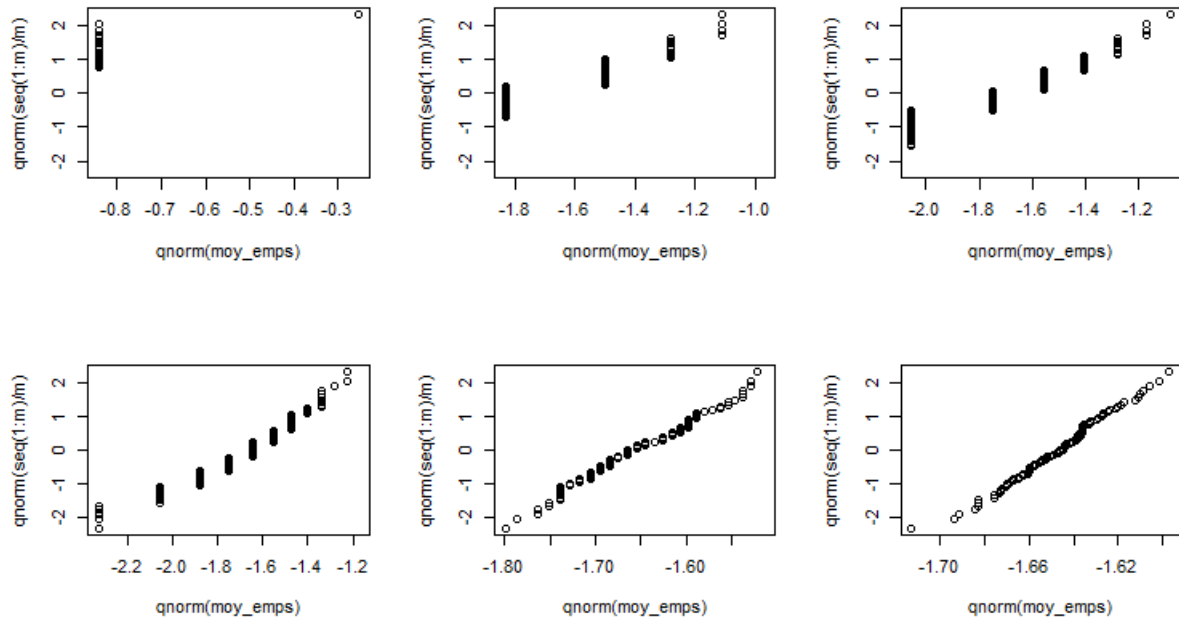
## Question 3

De la même manière que précédemment, on va faire la même simulation avec  $m = 100$  et en faisant varier  $n = 5, 30, 50, 100, 1000, 10000$ , sauf qu'on va tracer les histogrammes (de même largeur en l'occurrence) des moyennes empiriques pour chaque valeur de  $n$  ainsi que les graphes de probabilités associés pour la loi normale :





On remarque que la courbe associée à l'histogramme s'apparente à celle modélisant la fonction de densité de la loi normale à partir de  $n = 50$  environ.



De même, en traçant les graphes de probabilité  $(F_{\mathcal{N}}^{-1}(\bar{x}_i), F_{\mathcal{N}}^{-1}(\frac{i}{m}))$  avec  $\bar{x}_i \forall i = 1, \dots, n$  les moyennes empiriques de l'échantillon de données. On remarque de même que l'apparence d'une droite se forme à partir de  $n = 50$  ce qui confirme le fait qu'à partir de cette valeur de  $n = 50$  environ, la moyenne empirique des  $x_i$  (iid suivant une loi de Bernoulli) suit une loi normale. Le théorème centrale limite est ainsi illustrée.