

Compte-rendu en *Principes et méthodes statistiques*

Analyse de signaux oculométriques

Aurélien PEPIN, Léo DESBUREAUX, Julien LABOURÉ (Ensimag)

5 mai 2017

1 Analyse d'échantillons de loi binomiale négative

QUESTION 1. On suppose dans un premier temps r connu. L'estimateur obtenu par la méthode des moments est noté \tilde{p}_n .

Comme on suppose l'échantillon X_1, X_2, \dots, X_n indépendantes et de même loi binomiale négative $\mathcal{BN}(r, p)$, on a $E[X] = \frac{r}{p}$. Alors l'estimateur des moments est :

$$\tilde{p}_n = \frac{r}{\overline{X}_n}$$

où \overline{X}_n désigne la moyenne empirique de l'échantillon.

Une deuxième méthode d'estimation ponctuelle est l'estimation par maximum de vraisemblance. On note maintenant l'estimateur trouvé \hat{p}_n .

$$\mathcal{L}(p; x_1, \dots, x_n) = P(X_1 = x_1, \dots, X_n = x_n; p) = \prod_{i=1}^n P(X = x_i; p)$$

Plutôt que de dériver directement ce produit, on préfère maximiser le logarithme de la fonction de vraisemblance \mathcal{L} , c'est la **log-vraisemblance** :

$$\begin{aligned} \ln \mathcal{L}(p; x_1, \dots, x_n) &= \ln \prod_{i=1}^n P(X = x_i; p) = \sum_{i=1}^n \ln P(X = x_i; p) \\ &= \sum_{i=1}^n \ln \left(\binom{x_i-1}{r-1} (1-p)^{x_i-r} p^r \right) \end{aligned}$$

On cherche désormais à obtenir \hat{p}_n , valeur qui maximise cette log-vraisemblance. On dérive pour cela l'expression précédente :

$$\begin{aligned} \frac{\partial}{\partial p} \ln \mathcal{L}(p; x_1, \dots, x_n) &= \frac{\partial}{\partial p} \left(\sum_{i=1}^n \ln \binom{x_i-1}{r-1} \right) + \sum_{i=1}^n (x_i - r) \frac{\partial}{\partial p} (\ln(1-p)) + \sum_{i=1}^n r \frac{\partial}{\partial p} (\ln p) \\ &= 0 \qquad \qquad \qquad - \sum_{i=1}^n \left(\frac{x_i - r}{1-p} \right) \qquad \qquad \qquad + \sum_{i=1}^n \frac{r}{p} \end{aligned}$$

Cette expression s'annule sous les conditions suivantes :

$$\begin{aligned}
 -\sum_{i=1}^n \left(\frac{x_i - r}{1-p} \right) + \sum_{i=1}^n \frac{r}{p} = 0 &\iff \sum_{i=1}^n \frac{-px_i + pr + (1-p)r}{p(1-p)} = 0 \\
 &\iff \sum_{i=1}^n -px_i + r = 0 \\
 &\iff \sum_{i=1}^n \frac{r}{p} = \sum_{i=1}^n x_i \\
 &\iff n \frac{r}{p} = \sum_{i=1}^n x_i \\
 &\iff \frac{r}{p} = \bar{X}_n
 \end{aligned}$$

Finalement, on retrouve donc le résultat précédent :

$$\hat{p}_n = \frac{r}{\bar{X}_n} = \tilde{p}_n$$

Les cas aux limites où $p = 0$ ou $p = 1$ correspondent à des situations triviales où tous les X_i sont identiques, il n'y a aucune part d'aléatoire.

QUESTION 2. Pour une suite de variables aléatoires iid $\{X_n\}_{n \geq 1}$, le théorème central-limite exprime la convergence suivante :

$$Z_n = \sqrt{(n)} \frac{\bar{X}_n - E[X]}{\sigma(x)} \rightarrow \mathcal{N}(0, 1)$$

La suite X_1, \dots, X_n satisfait les conditions du théorème et on a les données suivantes :

$$\begin{aligned}
 - E[X] &= \frac{r}{p} \\
 - \sigma(X) &= \frac{\sqrt{(r(1-p))}}{p}
 \end{aligned}$$

Dans notre cas :

$$\begin{aligned}
 Z_n &= \sqrt{(n)} p \frac{\bar{X}_n - \frac{r}{p}}{\sqrt{(r(1-p))}} \\
 &= \frac{\sqrt{nr}}{\sqrt{1-p}} \left(\frac{p}{\hat{p}_n} - 1 \right) \rightarrow \mathcal{N}(0, 1)
 \end{aligned}$$

Sachant que $P(|Z_n| > u_\alpha) = \alpha$, on a :

$$\begin{aligned}
\frac{\sqrt{nr}}{\sqrt{1-p}} \left| \frac{p}{\hat{p}_n} - 1 \right| > u_\alpha &\iff \sqrt{rn} \left| p - \hat{p}_n \right| > u_\alpha \hat{p}_n \sqrt{1-p} \\
&\iff rn(p - \hat{p}_n)^2 > u_\alpha^2 \hat{p}_n^2 (1-p) \\
&\iff rn p^2 + (-2rn\hat{p}_n + u_\alpha^2 \hat{p}_n^2)p + \hat{p}_n^2 (rn - u_\alpha^2) > 0 \\
&\iff p^2 + \hat{p}_n(-2 + \frac{u_\alpha^2 \hat{p}_n}{rn})p + \hat{p}_n^2(1 - \frac{u_\alpha^2 \hat{p}_n}{rn}) > 0
\end{aligned}$$

Posons $\lambda = \frac{u_\alpha^2 \hat{p}_n}{rn}$, alors on a :

$$p^2 + p\hat{p}_n(-2 + \lambda\hat{p}_n) + \hat{p}_n^2(1 - \lambda) > 0$$

On obtient ainsi un polynôme en p de degré 2, dont on déduit le discriminant :


$$\Delta = (\hat{p}_n^2 \lambda)^2 + 1 + \frac{4(1 - \hat{p}_n)}{\lambda \hat{p}_n}$$

À partir du discriminant, on peut calculer les racines du polynôme qui sont les bornes de l'intervalle de confiance qu'on cherche à déterminer :

$$IC = \left[\hat{p}_n - \frac{1}{2} \lambda \hat{p}_n^2 \left(1 + \sqrt{1 + \frac{4(1 - \hat{p}_n)}{\lambda \hat{p}_n}} \right); \hat{p}_n - \frac{1}{2} \lambda \hat{p}_n^2 \left(1 - \sqrt{1 + \frac{4(1 - \hat{p}_n)}{\lambda \hat{p}_n}} \right) \right]$$

QUESTION 3. Pour tracer le graphe de probabilités de la loi géométrique (qui correspond au cas $r = 1$), on cherche en premier lieu des fonctions h, α, g, β telles que :

$$h[F(k)] = \alpha(p) g(k) + \beta(p)$$

 Se référer à : P1_Q3_Graphe_Probabilites.r

La fonction de répartition de la loi géométrique est :

$$\begin{aligned}
F_G(k) = 1 - (1-p)^k &\iff 1 - F_G(k) = (1-p)^k \\
&\iff \ln(1 - F_G(k)) = k \ln(1-p)
\end{aligned}$$

Par identification, on établit les correspondances suivantes :

- $h[F_G(k)] = \ln(1 - F_G(k))$
- $\alpha(p) = \ln(1 - p)$
- $g(k) = k$
- $\beta(p) = 0$

Le graphe de probabilités de $F_G(k)$ est donc le nuage de points :

$$(g(k_i^*); h(\frac{i}{n})) = (k_i^*; \ln(1 - \frac{i}{n})) \quad \forall i \in \llbracket 1; n-1 \rrbracket$$

INSERER ICI DE MIRIFIQUES GRAPHIQUES LE GROUPE 1 EST MEILLEUR QUE LE GROUPE 2

ON PREND ENSUITE LA PENTE DU GROUPE 1 ET ON L'INJECTE EN PARAMETRE D'UNE SIMU DE 10 000 DONNEES