
EXAMEN STATISTIQUE 1MF2E

Lundi 3 avril 2023

Partiel sans document (Une feuille A4 recto-verso autorisée)

Exercice 1 : Estimation (12 points)

On considère n observations x_1, \dots, x_n issues d'un échantillon (X_1, \dots, X_n) distribué suivant la même loi beta $B(1, \theta)$ de densité

$$f(x_i; \theta) = \begin{cases} \theta(1-x_i)^{\theta-1} & \text{si } x_i \in]0, 1[\\ 0 & \text{sinon} \end{cases}$$

avec $\theta > 0$.

1. (2pts) Montrer que l'estimateur du maximum de vraisemblance du paramètre θ noté $\hat{\theta}_{MV}$ en fonction des variables X_i est

$$\hat{\theta}_{MV} = -\frac{n}{\sum_{i=1}^n \ln(1-X_i)}.$$

On prendra soin de vérifier que cette vraisemblance admet un unique maximum global.

La vraisemblance de l'échantillon s'écrit

$$L(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta) = \theta^n \prod_{i=1}^n (1-x_i)^{\theta-1} \mathbb{I}_{]0,1[}(x_i)$$

où $\mathbb{I}_{]0,1[}$ est la fonction indicatrice sur l'intervalle $]0, 1[$. On sait qu'il est plus facile de travailler avec la log-vraisemblance définie par

$$\ln L(x_1, \dots, x_n; \theta) = n \ln \theta + (\theta - 1) \sum_{i=1}^n \ln(1-x_i).$$

Cette log-vraisemblance admet pour dérivée

$$\frac{\partial \ln L(x_1, \dots, x_n; \theta)}{\partial \theta} = \frac{n}{\theta} + \sum_{i=1}^n \ln(1-x_i).$$

Les variations de cette dérivée sont définies par

$$\begin{aligned} \frac{\partial \ln L(x_1, \dots, x_n; \theta)}{\partial \theta} \geq 0 &\Leftrightarrow \frac{n}{\theta} \geq -\sum_{i=1}^n \ln(1-x_i). \\ &\Leftrightarrow \theta \leq \frac{-n}{\sum_{i=1}^n \ln(1-x_i)}. \end{aligned}$$

Ceci permet de faire un tableau de variation qui indique que $\theta = \frac{-n}{\sum_{i=1}^n \ln(1-x_i)}$ est le maximum global unique de la vraisemblance, d'où

$$\hat{\theta}_{MV} = \frac{-n}{\sum_{i=1}^n \ln(1-X_i)}.$$

2. (2pts) Déterminer la loi de la variable aléatoire $Y_i = -\ln(1 - X_i)$ et en déduire sa fonction caractéristique (en s'aidant des tables). En déduire la loi de $T_n = -\sum_{i=1}^n \ln(1 - X_i)$. En utilisant le fait que si une variable Y suit une loi gamma $\mathcal{G}(\nu, \theta)$, son inverse $Z = \frac{1}{Y}$ suit une loi inverse gamma $\mathcal{IG}(\nu, \theta)$, montrer que $\hat{\theta}_{MV} = nZ_n$, où Z_n suit une loi inverse gamma $\mathcal{IG}(n, \theta)$. La variable aléatoire Y_i est à valeurs dans \mathbb{R}^+ . En faisant un changement de variables (et en faisant attention de ne pas oublier le Jacobien ;-)), on obtient la densité de Y_i

$$g(y_i; \theta) = \begin{cases} \theta e^{-\theta y_i} & \text{si } y_i > 0 \\ 0 & \text{sinon} \end{cases}$$

On reconnaît une loi gamma $\mathcal{G}(1, \theta)$ dont la fonction de répartition (donnée dans la table) est

$$\phi_{Y_i}(t) = E[e^{iY_i t}] = \frac{1}{1 - i\frac{t}{\theta}}.$$

La fonction caractéristique de T_n est le produit des fonctions caractéristiques des variables Y_i car ces variables sont indépendantes. On a donc

$$\phi_{T_n}(t) = E[e^{iT_n t}] = \frac{1}{(1 - i\frac{t}{\theta})^n}$$

qui est la fonction caractéristique d'une loi gamma $\mathcal{G}(n, \theta)$, donc $T_n \sim \mathcal{G}(n, \theta)$. D'après l'indication $Z_n = \frac{1}{T_n} = \frac{1}{\sum_{i=1}^n [-\ln(1 - X_i)]}$ suit une loi inverse gamma $\mathcal{IG}(n, \theta)$, d'où $\hat{\theta}_{MV} = nZ_n$, où Z_n suit une loi inverse gamma $\mathcal{IG}(n, \theta)$.

3. (2pts) Déterminer la moyenne de l'estimateur $\hat{\theta}_{MV}$. En déduire un estimateur sans biais et convergent du paramètre θ noté θ_n^* .

En utilisant les tables de lois et le fait que $\hat{\theta}_{MV} = nZ_n$, on obtient

$$E[\hat{\theta}_{MV}] = nE[Z_n] = n\frac{\theta}{n-1}.$$

$\hat{\theta}_{MV}$ est donc un estimateur asymptotiquement non biaisé de θ . On peut construire un estimateur non biaisé en multipliant $\hat{\theta}_{MV}$ par $\frac{n-1}{n}$. On obtient alors

$$\theta_n^* = (n-1)Z_n.$$

D'après la table, la variance de θ_n^* s'écrit

$$\text{var}[\theta_n^*] = (n-1)^2 \text{var}[Z_n] = (n-1)^2 \times \frac{\theta^2}{(n-1)^2(n-2)} = \frac{\theta^2}{n-2}.$$

Comme $E[\theta_n^*] - \theta = 0$ et que la variance de θ_n^* tend vers 0 lorsque $n \rightarrow \infty$, θ_n^* est un estimateur sans biais et convergent du paramètre θ .

4. (1pt) L'estimateur θ_n^* est-il l'estimateur efficace du paramètre θ ?

La dérivée seconde de la log-vraisemblance est

$$\frac{\partial^2 \ln L(x_1, \dots, x_n; \theta)}{\partial \theta^2} = \frac{-n}{\theta^2}.$$

d'où

$$E\left[-\frac{\partial^2 \ln L(X_1, \dots, X_n; \theta)}{\partial \theta^2}\right] = \frac{n}{\theta^2}.$$

On en déduit que la borne de Cramér-Rao pour un estimateur non-biaisé de θ est

$$\text{BCR} = \frac{-1}{E\left[\frac{\partial^2 \ln L(X_1, \dots, X_n; \theta)}{\partial \theta^2}\right]} = \frac{\theta^2}{n}.$$

Comme $\text{var}[\theta_n^*] > \frac{\theta^2}{n}$, θ_n^* n'est pas l'estimateur efficace du paramètre θ . On notera que cet estimateur est asymptotiquement efficace.

5. (5pts) On suppose maintenant qu'on dispose d'une information a priori sur le paramètre θ résumée dans la loi a priori de densité $p(\theta)$ définie par

$$p(\theta) = \lambda \exp(-\lambda\theta) \mathcal{I}_{\mathbb{R}^+}(\theta),$$

où $\mathcal{I}_{\mathbb{R}^+}$ est la fonction indicatrice sur \mathbb{R}^+ (on remarquera que cette loi est une loi gamma $\mathcal{G}(1, \lambda)$).

- Déterminer l'estimateur MAP de θ noté $\hat{\theta}_{\text{MAP}}$. Expliquer le comportement de $\hat{\theta}_{\text{MAP}}$ et de $\hat{\theta}_{\text{MV}}$ lorsque n tend vers $+\infty$.

La loi a posteriori du paramètre θ vérifie

$$\begin{aligned} p(\theta|x_1, \dots, x_n) &\propto p(x_1, \dots, x_n|\theta)\pi(\theta) \\ &\propto \theta^n \exp\left\{-\theta \left[\lambda - \sum_{i=1}^n \ln(1-x_i)\right]\right\} \mathcal{I}_{]0,+\infty[}(\theta), \end{aligned}$$

qui a la même forme que la log-vraisemblance si on remplace $\sum_{i=1}^n \ln(1-x_i)$ par $\lambda - \sum_{i=1}^n \ln(1-x_i)$. L'estimateur du maximum a posteriori du paramètre θ s'obtient donc à partir de l'estimateur du maximum de vraisemblance en faisant cette même transformation, soit

$$\hat{\theta}_{\text{MAP}} = \frac{n}{\lambda - \sum_{i=1}^n \ln(1-X_i)}.$$

Remarque : l'estimateur du maximum a posteriori du paramètre θ s'écrit

$$\hat{\theta}_{\text{MAP}} = \frac{1}{\frac{\lambda}{n} + \frac{1}{\hat{\theta}_{\text{MV}}}}$$

On peut donc remarquer que $\hat{\theta}_{\text{MAP}}$ et $\hat{\theta}_{\text{MV}}$ se comportent de manière similaire lorsque $n \rightarrow +\infty$. Lorsqu'on a beaucoup de données, on a tendance à oublier l'information a priori et à faire confiance à l'estimateur du maximum de vraisemblance qui résume l'information contenue dans les données.

- Montrer que la loi a posteriori de $\theta|x_1, \dots, x_n$ est une loi gamma dont on déterminera les paramètres. En déduire l'estimateur de la moyenne a posteriori du paramètre θ noté $\hat{\theta}_{\text{MMSE}}$. La loi a posteriori du paramètre θ vérifie

$$p(\theta|x_1, \dots, x_n) \propto \theta^n \exp\left\{-\theta \left[\lambda - \sum_{i=1}^n \ln(1-x_i)\right]\right\} \mathcal{I}_{]0,+\infty[}(\theta).$$

qui est une loi gamma $\mathcal{G}(n+1, \lambda - \sum_{i=1}^n \ln(1-x_i))$. On en déduit

$$\hat{\theta}_{\text{MMSE}} = E[\theta|X_1, \dots, X_n] = \frac{n+1}{\lambda - \sum_{i=1}^n \ln(1-X_i)}.$$

Exercice 2 : Test Statistique (8 points)

On considère un ensemble de variables aléatoires (X_1, \dots, X_n) indépendantes et de densités

$$f(x_i; \theta) = \begin{cases} \theta(1 - x_i)^{\theta-1} & \text{si } x_i \in]0, 1[\\ 0 & \text{sinon} \end{cases}$$

avec $\theta > 0$. On désire tester les deux hypothèses

$$H_0 : \theta = \theta_0, \quad H_1 : \theta = \theta_1 > \theta_0.$$

1. (2pts) À l'aide du théorème de Neyman-Pearson, montrer que la statistique du test le plus puissant est $T_n = -\sum_{i=1}^n \ln(1 - X_i)$ et indiquer la région critique de ce test.

Le test de Neyman Pearson est défini par

$$\text{Rejet de } H_0 \text{ si } \frac{L(x_1, \dots, x_n; \theta_1)}{L(x_1, \dots, x_n; \theta_0)} > S_{1,\alpha}$$

où $S_{1,\alpha}$ est un seuil dépendant du risque de première espèce α . Mais

$$\begin{aligned} \frac{L(x_1, \dots, x_n; \theta_1)}{L(x_1, \dots, x_n; \theta_0)} > S_{1,\alpha} &\Leftrightarrow \ln \frac{L(x_1, \dots, x_n; \theta_1)}{L(x_1, \dots, x_n; \theta_0)} > S_{2,\alpha} \\ &\Leftrightarrow n \ln \left(\frac{\theta_1}{\theta_0} \right) + (\theta_1 - \theta_0) \sum_{i=1}^n \ln(1 - x_i) > S_{2,\alpha}. \end{aligned}$$

Comme $\theta_1 > \theta_0$, on a la règle de décision

$$\text{Rejet de } H_0 \text{ si } T_n < S_\alpha.$$

La région critique du test est l'ensemble des vecteurs $(x_1, \dots, x_n) \in]0, 1[^n$ tels que $T_n < S_\alpha$ et la statistique de test est T_n .

2. (1pt) On admet que si X_i suit une loi de densité $f(x_i; \theta)$, alors $Y_i = -\ln(1 - X_i)$ suit une loi gamma $\mathcal{G}(1, \theta)$. En utilisant le théorème central limite, déterminer la loi asymptotique de T_n sous les deux hypothèses.

La loi gamma $\mathcal{G}(1, \theta)$ est de moyenne $\frac{1}{\theta}$ et de variance $\frac{1}{\theta^2}$. Donc

$$E[T_n] = \frac{n}{\theta} \text{ et } \text{var}[T_n] = \frac{n}{\theta^2}.$$

Comme T_n peut s'écrire comme une somme de variables indépendantes, l'application du théorème central limite permet alors d'obtenir les résultats suivants

$$\text{Sous } H_0 : T_n \approx \mathcal{N} \left(\frac{n}{\theta_0}, \frac{n}{\theta_0^2} \right),$$

$$\text{Sous } H_1 : T_n \approx \mathcal{N} \left(\frac{n}{\theta_1}, \frac{n}{\theta_1^2} \right),$$

où \approx signifie "de loi approchée (pour n grand)".

3. (2pts) En utilisant la loi asymptotique de la question précédente, déterminer les risques de première et seconde espèce α et β du test en fonction des paramètres θ_0 et θ_1 et de la fonction de répartition d'une loi normale $\mathcal{N}(0, 1)$ notée F .

Le risque α est défini par

$$\alpha = P[\text{Rejeter } H_0 | H_0 \text{ vraie}] = P \left[T_n < S_\alpha | T_n \sim \mathcal{N} \left(\frac{n}{\theta_0}, \frac{n}{\theta_0^2} \right) \right],$$

soit

$$\alpha = F \left[\frac{S_\alpha - \frac{n}{\theta_0}}{\sqrt{\frac{n}{\theta_0^2}}} \right].$$

De même

$$\beta = P[\text{Rejeter } H_1 | H_1 \text{ vraie}] = P \left[T_n \geq S_\alpha | T_n \sim \mathcal{N} \left(\frac{n}{\theta_1}, \frac{n}{\theta_1^2} \right) \right],$$

soit

$$\beta = 1 - F \left[\frac{S_\alpha - \frac{n}{\theta_1}}{\sqrt{\frac{n}{\theta_1^2}}} \right].$$

4. (2pts) Déterminer les courbes COR associées à ce test montrer qu'elles ne dépendent que de n et de $\frac{\theta_1}{\theta_0}$. Tracer l'allure de ces courbes COR pour différentes valeurs de n dans une première figure et pour différentes valeurs de $\frac{\theta_1}{\theta_0}$ dans une seconde figure.

Les courbes COR expriment $\pi = 1 - \beta$ en fonction de α . On a donc

$$\pi = F \left[\frac{S_\alpha - \frac{n}{\theta_1}}{\sqrt{\frac{n}{\theta_1^2}}} \right]$$

avec $S_\alpha = F^{-1}(\alpha) \sqrt{\frac{n}{\theta_0^2}} + \frac{n}{\theta_0}$, soit

$$\pi = F \left[F^{-1}(\alpha) \frac{\theta_1}{\theta_0} + \sqrt{n} \left(\frac{\theta_1}{\theta_0} - 1 \right) \right].$$

On observe donc que π ne dépend que de n et de $\frac{\theta_1}{\theta_0}$. De plus, comme F est une fonction croissante (c'est une fonction de répartition), π est une fonction croissante de n , ce qui est habituel car plus le nombre d'observations est grand, meilleure est la performance du test. La puissance du test est aussi une fonction croissante de $\frac{\theta_1}{\theta_0}$, ce qui donne des courbes COR similaires à celles représentées ci-dessous

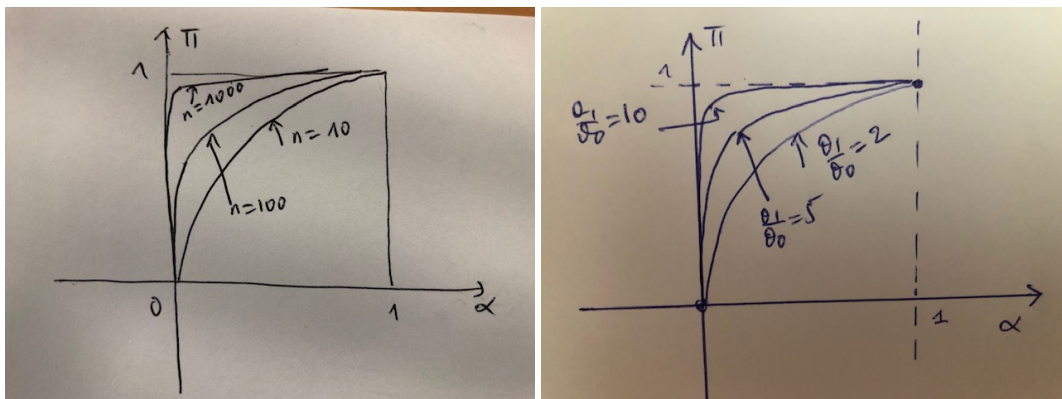


Figure 1: Allure des courbes COR pour différentes valeurs de n et de $\frac{\theta_1}{\theta_0}$.

5. (3pts) On désire vérifier que l'ensemble des observations (x_1, \dots, x_n) suit une loi beta $B(1, \theta)$ de paramètre $\theta = \frac{1}{2}$ à l'aide d'un test du χ^2 . Déterminer la fonction de répartition de cette loi et en déduire que l'intervalle $]0, 1[$ est la réunion de trois intervalles équiprobables pour la loi $B(1, \frac{1}{2})$ que l'on précisera. On compte le nombre d'observations x_i appartenant à ces trois intervalles et on trouve $K_1 = 13$, $K_2 = 8$ et $K_3 = 9$. Quelle est la valeur de la statistique du test du χ^2 ? Exprimer

le seuil de ce test noté S_α en fonction du risque α et de l'inverse de la fonction de répartition d'une loi du χ^2 dont on précisera le nombre de degrés de liberté. On donne $S_{0.05} = 5.991$. Qu'en conclut-on ?

La fonction de répartition d'une loi $B(1, \theta)$ s'écrit

$$F_\theta(x) = \int_{-\infty}^x \theta(1-u)^{\theta-1} du = \begin{cases} 0 & \text{si } x \leq 0 \\ 1 - (1-x)^\theta & \text{si } x \in]0, 1[\\ 1 & \text{si } x \geq 1 \end{cases}$$

Pour $\theta = \frac{1}{2}$, on obtient donc

$$F_{\frac{1}{2}}(x) = \begin{cases} 0 & \text{si } x \leq 0 \\ 1 - \sqrt{1-x} & \text{si } x \in]0, 1[\\ 1 & \text{si } x \geq 1 \end{cases}$$

Le premier intervalle $]0, a[$ est tel que

$$F(a) = \frac{1}{3} \Leftrightarrow 1 - \sqrt{1-x} = \frac{1}{3} \Leftrightarrow x = \frac{5}{9}$$

Le deuxième intervalle $[a, b]$ est tel que

$$F(b) = \frac{2}{3} \Leftrightarrow 1 - \sqrt{1-x} = \frac{2}{3} \Leftrightarrow x = \frac{8}{9}$$

La statistique du test du χ^2 est

$$\phi = \sum_{i=1}^3 \frac{(K_i - np_i)^2}{np_i} = \frac{9}{10} + \frac{4}{10} + \frac{1}{10} = \frac{14}{10} = 1.4.$$

On sait que sous l'hypothèse H_0 , ϕ suit une loi du χ^2 à $K - 1$ degrés de liberté, donc

$$\phi \sim \chi_2^2.$$

Comme $\phi = 1.4 < S_{0.05} = 5.991$, on accepte l'hypothèse H_0 avec le risque $\alpha = 0.05$ donc on en déduit que l'échantillon est de loi $B(1, \frac{1}{2})$ avec le risque $\alpha = 0.05$.

LOIS DE PROBABILITÉ CONTINUES

m : moyenne σ^2 : variance F. C. : fonction caractéristique

LOI	Densité de probabilité	m	σ^2	F. C.
Uniforme	$f(x) = \frac{1}{b-a}$ $x \in]a, b[$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{itb} - e^{ita}}{it(b-a)}$
Gamma $\mathcal{G}(\nu, \theta)$	$f(x) = \frac{\theta^\nu}{\Gamma(\nu)} e^{-\theta x} x^{\nu-1}$ $\theta > 0, \nu > 0$ $x \geq 0$ avec $\Gamma(n+1) = n! \forall n \in \mathbb{N}$	$\frac{\nu}{\theta}$	$\frac{\nu}{\theta^2}$	$\frac{1}{(1 - i\frac{t}{\theta})^\nu}$
Inverse gamma $\mathcal{IG}(\nu, \theta)$	$f(x) = \frac{\theta^\nu}{\Gamma(\nu)} e^{-\frac{\theta}{x}} \frac{1}{x^{\nu+1}}$ $\theta > 0, \nu > 0$ $x \geq 0$ avec $\Gamma(n+1) = n! \forall n \in \mathbb{N}$	$\frac{\theta}{\nu-1}$ si $\nu > 1$	$\frac{\theta^2}{(\nu-1)^2(\nu-2)}$ si $\nu > 2$	(*)
Première loi de Laplace	$f(x) = \frac{1}{2} e^{- x }, \quad x \in \mathbb{R}$	0	2	$\frac{1}{1+t^2}$
Normale univariée $\mathcal{N}(m, \sigma^2)$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}, \quad x \in \mathbb{R}$	m	σ^2	$e^{imt - \frac{\sigma^2 t^2}{2}}$
Normale multivariée $\mathcal{N}_p(\mathbf{m}, \Sigma)$	$f(\mathbf{x}) = K e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \Sigma^{-1}(\mathbf{x}-\mathbf{m})}$ $K = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}}$ $\mathbf{x} = (x_1, \dots, x_p)^T \in \mathbb{R}^p$	\mathbf{m}	Σ	$e^{i\mathbf{u}^T \mathbf{m} - \frac{1}{2} \mathbf{u}^T \Sigma \mathbf{u}}$
Khi2 χ_ν^2 $\Gamma(\frac{1}{2}, \frac{\nu}{2})$	$f(x) = k e^{-\frac{x}{2}} x^{\frac{\nu}{2}-1}$ $k = \frac{1}{2^{\frac{\nu}{2}} \Gamma(\frac{\nu}{2})}$ $\nu \in \mathbb{N}^*, x \geq 0$	ν	2ν	$\frac{1}{(1-2it)^{\frac{\nu}{2}}}$
Cauchy $c_{\lambda, \alpha}$	$f(x) = \frac{1}{\pi \lambda \left(1 + \left(\frac{x-\alpha}{\lambda}\right)^2\right)}$ $\lambda > 0, \alpha \in \mathbb{R}$	(-)	(-)	$e^{i\alpha t - \lambda t }$
Beta $B(a, b)$	$f(x) = k x^{a-1} (1-x)^{b-1}$ $k = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$ $a > 0, b > 0$ $x \in]0, 1[$ avec $\Gamma(n+1) = n! \forall n \in \mathbb{N}$	$\frac{a}{a+b}$	$\frac{ab}{(a+b)^2(a+b+1)}$	(*)

LOIS DE PROBABILITÉ DISCRÈTES

m : moyenne σ^2 : variance **F. C.** : fonction caractéristique

$p_k = P[X = k]$ $p_{1,\dots,m} = P[X_1 = k_1, \dots, X_m = k_m]$

LOI	Probabilités	m	σ^2	F. C.
Uniforme	$p_k = \frac{1}{n}$ $k \in \{1, \dots, n\}$	$\frac{n+1}{2}$	$\frac{n^2-1}{12}$	$\frac{e^{it}(1 - e^{itn})}{n(1 - e^{it})}$
Bernoulli	$p_1 = P[X = 1] = p$ $p_0 = P[X = 0] = q$ $p \in [0, 1]$ $q = 1 - p$	p	pq	$pe^{it} + q$
Binomiale $B(n, p)$	$p_k = C_n^k p^k q^{n-k}$ $p \in [0, 1]$ $q = 1 - p$ $k \in \{0, 1, \dots, n\}$	np	npq	$(pe^{it} + q)^n$
Binomiale négative	$p_k = C_{n+k-1}^{n-1} p^n q^k$ $p \in [0, 1]$ $q = 1 - p$ $k \in \mathbb{N}$	$n \frac{q}{p}$	$n \frac{q}{p^2}$	$\left(\frac{p}{1 - qe^{it}}\right)^n$
Multinomiale	$p_{1,\dots,m} = \frac{n!}{k_1! \dots k_m!} p_1^{k_1} \dots p_m^{k_m}$ $p_j \in [0, 1]$ $q_j = 1 - p_j$ $k_j \in \{0, 1, \dots, n\}$ $\sum_{j=1}^m k_j = n$ $\sum_{j=1}^m p_j = 1$	np_j	Variance : $np_j q_j$ Covariance : $-np_j p_k$	$\left(\sum_{j=1}^m p_j e^{it}\right)^n$
Poisson $P(\lambda)$	$p_k = e^{-\lambda} \frac{\lambda^k}{k!}$ $\lambda > 0$ $k \in \mathbb{N}$	λ	λ	$\exp[\lambda(e^{it} - 1)]$
Géométrique	$p_k = pq^{k-1}$ $p \in [0, 1]$ $q = 1 - p$ $k \in \mathbb{N}^*$	$\frac{1}{p}$	$\frac{q}{p^2}$	$\frac{pe^{it}}{1 - qe^{it}}$