

## Contribution à l'évaluation d'un patrimoine iconographique en vue d'applications aux pays d'Afrique subsaharienne Mesure des disparités liées à la corrélation des caractères temporel et régional

*par Stéphane Richemond\**

Nous avons souvent été confrontés aux disparités des patrimoines iconographiques. Ce fut le cas pour les expositions que I&M avait préparées sur les villes de Bamako, Lomé, Ouagadougou... ou sur le Vietnam. Nous avons pensé qu'il n'était pas suffisant de faire le constat de l'absence d'images sur tel thème, tel quartier ou telle région, telle période, mais qu'il serait souhaitable de quantifier ces carences de façon objective à l'aide d'indices qui donneraient une mesure des disparités permettant de caractériser un patrimoine et de comparer des patrimoines différents. Cette idée, qui semblait ne pas avoir été explorée, nous a conduit à établir les formules présentées dans cet article et dans les deux précédentes livraisons de ce *Bulletin*<sup>1</sup>. Nous avons vu comment mesurer les disparités temporelles et régionales du patrimoine iconographique d'un pays<sup>2</sup>. Après nous être intéressés à la distribution d'une population d'images selon les deux caractères, la période *P* et le lieu (ou la région) *R*, ceux-ci étant pris isolément, nous présentons ici une mesure des disparités dues à leur corrélation.

### *I - Rappel sur la mesure des disparités temporelles d'une population d'images*

Une population d'images couvre toujours de façon inégale les périodes considérées. Supposons que l'intervalle de temps étudié soit la réunion de *m* périodes successives. Soit *N* le nombre total d'images toutes périodes confondues. Notons *n<sub>i</sub>* le nombre d'images appartenant à la période *P<sub>i</sub>* et *f<sub>i</sub>* la proportion *n<sub>i</sub>/N* d'images appartenant à cette période. Selon le caractère *P*, la distribution (dite marginale) d'images peut être décrite dans le tableau suivant :

<i>P<sub>i</sub></i>	<i>P<sub>1</sub></i>	<i>P<sub>2</sub></i>		<i>P<sub>i</sub></i>		<i>P<sub>m</sub></i>	Total
<i>f<sub>i</sub></i>	<i>f<sub>1</sub></i>	<i>f<sub>2</sub></i>		<i>f<sub>i</sub></i>	.	<i>f<sub>m</sub></i>	1

Supposons les périodes *P<sub>i</sub>* d'égale longueur, alors on peut considérer que la population d'images serait harmonieusement distribuée si chaque période contenait le même nombre d'images *N/m* soit la même proportion *p<sub>i</sub>=1/m* d'images. Une telle distribution idéale peut être représentée dans le tableau suivant :

<i>P<sub>i</sub></i>	<i>P<sub>1</sub></i>	<i>P<sub>2</sub></i>		<i>P<sub>i</sub></i>			<i>P<sub>m</sub></i>	Total
<i>p<sub>i</sub></i>	1/m	1/m		1/m			1/m	1

Nous pouvons mesurer l'écart entre la distribution d'images et la distribution idéale de référence par l'indicateur suivant que nous avons appelé **indice de disparité temporelle**<sup>3</sup> :

\* IRHiS (Université de Lille) - srichemond@hotmail.com

<sup>1</sup> - Stéphane Richemond, "Contribution à l'évaluation d'un patrimoine iconographique en vue d'applications aux pays d'Afrique subsaharienne - Mesure des disparités temporelles", *Bulletin n°53*, Images & Mémoires, Été 2017.

- Stéphane Richemond, "Contribution à l'évaluation d'un patrimoine iconographique en vue d'applications aux pays d'Afrique subsaharienne - Mesure des disparités régionales", *Bulletin n°54*, Images & Mémoires, Automne 2017.

<sup>2</sup> Nous avons défini un tel patrimoine par l'ensemble des images (dessins, peintures, photographies,...) identifiables comme appartenant au pays considéré, où qu'elles soient localisées, même si elles ne sont pas à la disposition du public.

<sup>3</sup> Il s'agit de la distance euclidienne définie par le théorème de Pythagore généralisé. En divisant *E<sub>p</sub>* par sa valeur maximum prise lorsque toutes les images sont concentrées sur une seule période, on obtient un indice compris entre 0 et 1.

Nous avons vu précédemment que nous pouvions aussi utiliser la distance du khi<sup>2</sup> définie par :

$$\sqrt{\sum_{i=1}^m \left( \frac{f_i - p_i}{\sqrt{p_i}} \right)^2}$$

$$E_p = \sqrt{\sum_{i=1}^m (f_i - \frac{1}{m})^2} \text{ où on a } \sum_{i=1}^m n_i = N \text{ et } \sum_{i=1}^m f_i = 1$$

**II - Rappel sur la mesure des disparités régionales d'une population d'images**

Il nous a semblé naturel de considérer une distribution uniforme comme référence pour la mesure des disparités temporelles ce qui n'est pas justifié pour la mesure des disparités régionales car les régions ont des superficies et des populations différentes. Nous avons pensé pertinent de considérer des proportions de référence proportionnelles aux superficies des régions, ou encore à leurs populations. Supposons que l'étendue E du pays soit la réunion de p régions R<sub>j</sub>. Notons n<sub>j</sub> le nombre d'images concernant la région R<sub>j</sub> et f<sub>j</sub> = n<sub>j</sub>/N la proportion correspondante. La distribution (dite marginale) d'images est alors la suivante :

R <sub>1</sub>	R <sub>1</sub>	R <sub>2</sub>		R <sub>j</sub>		R <sub>p</sub>	Total
f <sub>1</sub>	f <sub>1</sub>	f <sub>2</sub>		f <sub>j</sub>	.	f <sub>p</sub>	1

Notons S la superficie de l'ensemble des régions considérées, et S<sub>j</sub> la superficie de la région R<sub>j</sub>, de même A la population totale et A<sub>j</sub> celle de la région R<sub>j</sub>. Prenons pour proportion de référence p<sub>j</sub> = S<sub>j</sub>/S. Nous pouvons mesurer l'écart entre la distribution d'images et la distribution de référence par l'expression suivante que nous avons appelée **indice de disparité régionale relatif aux superficies** :

$$E_{R/S} = \sqrt{\sum_{j=1}^p (f_j - \frac{S_j}{S})^2}$$

Définissons maintenant la proportion de référence par p<sub>j</sub> = A<sub>j</sub>/A. Nous avons mesuré l'écart entre la distribution d'images et la distribution de référence par l'expression suivante que nous avons appelé **indice de disparité régionale relatif aux populations** :

$$E_{R/A} = \sqrt{\sum_{j=1}^p (f_j - \frac{A_j}{A})^2}$$

Par ailleurs, comme annoncé précédemment, on adoptera plus volontiers la formule de la distance du khi<sup>2</sup> pour chacun de ces indices. Pour prendre en compte à la fois les superficies et les proportions pour chaque région, nous pourrions prendre des proportions de référence du type :

$$p_j = \alpha \frac{S_j}{S} + (1 - \alpha) \frac{A_j}{A} \text{ où } \alpha \in [0,1]$$

**III - Tableau de contingence**

On s'intéresse maintenant à l'étude une population d'images sous l'angle simultané des deux caractères la période et la région. On notera n<sub>ij</sub> (respectivement f<sub>ij</sub>) le nombre (respectivement la proportion) d'images qui possèdent à la fois la modalité P<sub>i</sub> du caractère P et la modalité R<sub>j</sub> du caractère R. On a ainsi :

$$f_{i.} = \sum_{j=1}^p f_{ij}; \text{ de même } f_{.j} = \sum_{i=1}^m f_{ij} \text{ et } \sum_{i=1}^m \sum_{j=1}^p f_{ij} = 1$$

où le point en indice signifie que l'on a effectué une sommation sur l'indice remplacé. On peut ranger les effectifs ou les proportions dans un tableau à double entrée :

	R <sub>1</sub>	R <sub>2</sub>		R <sub>j</sub>		R <sub>p</sub>	f <sub>i.</sub>
P <sub>1</sub>	f <sub>11</sub>	f <sub>12</sub>		f <sub>1j</sub>		f <sub>1p</sub>	f <sub>1.</sub>
P <sub>2</sub>	f <sub>21</sub>	f <sub>22</sub>		f <sub>2j</sub>		f <sub>2p</sub>	f <sub>2.</sub>
P <sub>i</sub>	f <sub>i1</sub>	f <sub>i2</sub>		f <sub>ij</sub>		f <sub>ip</sub>	f <sub>i.</sub>
P <sub>m</sub>	f <sub>m.</sub>	f <sub>m.</sub>		f <sub>mi</sub>		f <sub>m.</sub>	f <sub>m.</sub>
f <sub>.j</sub>	f <sub>.1</sub>	f <sub>.2</sub>		f <sub>.j</sub>		f <sub>.p</sub>	1

Ainsi, on a porté à l'intersection de la ième ligne et de la jème colonne la proportion  $f_{ij}$  des images qui appartiennent à la fois à la période  $P_i$  et à la région  $R_j$ . Nous avons déjà insisté sur le grand intérêt qu'il y a à remplir ce tableau qui met en évidence les disparités des caractères en un coup d'œil.

**IV - Distributions conditionnelles**

On peut aussi étudier la population d'images sous l'angle d'un caractère, les modalités de l'autre étant fixée. Par exemple, on peut étudier le caractère  $P$  dans le cas où la modalité  $R_j$  de  $R$  est fixée. On obtient le tableau :

P <sub>i</sub>	P <sub>1</sub>	P <sub>2</sub>		P <sub>i</sub>		P <sub>m</sub>	Total
f <sub>i</sub> <sup>j</sup>	f <sub>1</sub> <sup>j</sup>	f <sub>2</sub> <sup>j</sup>		f <sub>i</sub> <sup>j</sup>	.	f <sub>m</sub> <sup>j</sup>	f <sup>j</sup>

où on a noté  $f_i^j = n_{ij}/n_j$  ce qui se lit « f<sub>i</sub> si j ».

On peut donc étudier, sous l'angle du caractère  $P$ , p distributions à j fixé. De même, on peut étudier, sous l'angle du caractère  $R$ , m distributions à i fixé.

Il est clair que caractère  $P$  est indépendant du caractère  $R$  si la proportion d'individus possédant la modalité  $P_i$  de  $P$  reste constante pas lorsque j varie. En supposant que  $P$  et  $R$  soient les seuls caractères, on a alors :

$$f_i^j = \text{cte } \forall j \in \{1, 2, \dots, p\} \text{ soit } f_i^j = f_i.$$

On montre aisément que l'indépendance est une notion réciproque. En effet :

$$f_{ij} = \frac{n_{ij}}{N} = \frac{n_{ij}}{n_j} \cdot \frac{n_j}{N} = f_i^j \cdot f_j = f_i \cdot f_j \text{ et } f_{ij} = f_i^i \cdot f_i. \text{ On en déduit : } f_i^j = f_j.$$

On observe alors que les lignes sont proportionnelles entre elles, de même que les colonnes le sont aussi.

**V - Disparités des distributions conditionnelles**

Nous avons déjà étudié les disparités des distributions marginales de  $P$  et  $R$ . Les indices  $E_R$  et  $E_P$  ne sont cependant pas suffisants pour caractériser toutes les disparités d'une population. En effet, imaginons le cas extrême où :  $f_i = f_j \forall i \in \{1, 2, \dots, m\}, \forall j \in \{1, 2, \dots, p\}, f_{ij} = 0$  si  $i \neq j$ , soit, par exemple, à la situation du tableau suivant :

	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	R <sub>4</sub>	n <sub>i.</sub>
P <sub>1</sub>	0,25	0	0	0	0,25
P <sub>2</sub>	0	0,25	0	0	0,25
P <sub>3</sub>	0	0	0,25	0	0,25
P <sub>4</sub>	0	0	0	0,25	0,25
n <sub>.j</sub>	0,25	0,25	0,25	0,25	1

Alors que les distributions marginales de  $P$  et  $R$  montrent des disparités nulles ( $E_R = E_P = 0$ ), nous avons une disparité maximum des distributions conditionnelles. En effet, toutes les images de la période  $P_1$  sont concentrées sur la région  $R_1$ , et ainsi de suite. Ce cas correspond au maximum de corrélation (liaison fonctionnelle réciproque des deux caractères  $P$  et  $R$ ).

La disparité des distributions conditionnelles peut être mesurée par la distance entre le tableau de corrélation et le tableau correspondant à l'indépendance des caractères  $P$  et  $R$  qui est celui de la population de référence idéale. Le tableau correspondant à l'indépendance des caractères se construit aisément à partir de la connaissance des deux distributions marginales grâce à la relation :

$$f_{ij} = f_i \cdot f_j \quad \forall i \in \{1, 2, \dots, m\}, \forall j \in \{1, 2, \dots, p\} \quad \text{qui s'écrit encore } n_{ij} = n_i \cdot f_j.$$

ce qui donnerait ici le tableau de référence (comprenant des proportions égales dans ce cas particulier) :

	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	R <sub>4</sub>	f <sub>i</sub>
P <sub>1</sub>	0,0625	0,0625	0,0625	0,0625	0,25
P <sub>2</sub>	0,0625	0,0625	0,0625	0,0625	0,25
P <sub>3</sub>	0,0625	0,0625	0,0625	0,0625	0,25
P <sub>4</sub>	0,0625	0,0625	0,0625	0,0625	0,25
f <sub>j</sub>	0,25	0,25	0,25	0,25	1

D'une façon plus générale, nous définirons **l'indice de disparité des distributions conditionnelles** par :

$$E_{PR} = \sqrt{\sum_{i=1}^m \sum_{j=1}^p (f_{ij} - f_i \cdot f_j)^2}$$

qui correspond à une mesure de la distance entre les deux tableaux précédents.

En pratique, les proportions du patrimoine iconographique ne sont pas connues et  $E_{PR}$  n'est donc pas calculable par la formule ci-dessus. Nous utiliserons cette formule en remplaçant les proportions par celles déterminées sur un échantillon aléatoire d'images, comme nous avons procédé lors des études précédentes concernant les disparités temporelles et régionales. Il serait souhaitable que l'échantillon soit le plus grand possible, mais il est surtout important que l'échantillon, qui est généralement une collection, n'ait pas été constitué avec des images appartenant à une période, une région ou un thème privilégié.

La quantité obtenue à partir des proportions  $f_{ij}^*$  déterminées sur l'échantillon aléatoire :

$$E_{PR}^* = \sqrt{\sum_{i=1}^m \sum_{j=1}^p (f_{ij}^* - f_i^* \cdot f_j^*)^2}$$

est une variable qui est un estimateur asymptotiquement sans biais<sup>4</sup> de  $E_{PR}$ , condition nécessaire mais non suffisante à sa convergence en probabilité vers  $E_{PR}$  lorsque la taille de l'échantillon  $n$  croît.

### Conclusion

Cette étude ouvre la voie à des développements plus importants. En particulier, nous n'avons pas abordé la mesure des disparités thématiques pour lesquelles la définition d'une population de référence est délicate. Une réflexion sur la détermination d'un indice plus synthétique pourrait être menée. La question d'établir un échantillon aléatoire à partir de collections d'images n'est pas non plus un problème simple.

<sup>4</sup>

$$\text{Ecrivons } \sum_{i=1}^m \sum_{j=1}^p (f_{ij}^* - f_i^* \cdot f_j^*)^2 \text{ sous la forme : } \sum_{i=1}^m \sum_{j=1}^p ((f_{ij}^* - f_{ij}) + (f_{ij} - f_i \cdot f_j) + (f_{ij} \cdot f_j - f_i^* \cdot f_j^*))^2$$

que nous développerons, obtenant ainsi la somme de six termes dont l'espérance mathématique tend vers  $E_{PR}$  lorsque  $n$  croît. Pour l'établir, on aura recours à :

$$E(f_i^* \cdot f_j^* - f_i \cdot f_j) = \text{Cov}(f_i^*, f_j^*) \leq \sigma(f_i^*) \sigma(f_j^*) = \frac{1}{n} \sqrt{(1 - f_i)(1 - f_j) f_i \cdot f_j} \quad \text{qui tend vers 0 avec } n.$$

Ici Cov représente la covariance et  $\sigma$  représente l'écart-type.