

ON SOME PROPERTIES ON CHI-SQUARE DISTANCE

Khaled Moussa Aludaat

Department of Statistics, Yarmouk University,
Yarmouk, Jordan.

ABSTRACT

In this paper we have described the properties of the distance between two multi dimensional points, which is called Chi-square distance. It is a weighted Euclidean distance.

INTRODUCTION

A large number of multivariate problem can be viewed in terms of "distances" between two observations, or between samples of observations, or between populations of observations. Many distance measures have been proposed and used in multivariate analysis. Here we introduce the distance used in factor analysis of correspondence (method of a multivariate analysis proposed by J.P. Benzercri in (1962)).

This distance has many properties which deserves our attention.

Let K_{IJ} be a contingency table of positive numbers $k(i,j)$ with I rows and J columns:

$$K_{IJ} = \{k(i,j) \mid i \in I, j \in J\}$$

We define the marginals k_i , k_j and the total K over K_{IJ} ; also the profiles f_j^i and f_i^j of the rows and the columns of K_{IJ} . The definitions are given below:

$$\forall i \in I, k(i) = \sum_j k(i,j), K_i = \{k(i) \mid i \in I\}$$

$$\forall j \in J, k(j) = \sum_i k(i,j), K_j = \{k(j) \mid j \in J\}$$

$$K = \sum_i \sum_j k(i,j).$$

$$f_{ij} = \frac{k(i,j)}{K}, f_{IJ} = \{f_{ij} \mid i \in I, j \in J\}$$

$$f_i = \frac{k(i)}{K}, f_I = \{f_i \mid i \in I\}$$

$$f_j = \frac{k(j)}{K}, f_J = \{f_j \mid j \in J\}$$

$$f_j^i = \frac{k(i,j)}{K(i)} = \frac{f_{ij}}{f_i}, f_i^j = \{f_j^i \mid j \in J\}$$

$$f_j^i = \frac{k(i,j)}{K(j)} = \frac{f_{ij}}{f_j}, f_i^j = \{f_j^i \mid i \in I\}$$

The class of the profiles of the rows is defined as:

$$N_I = \{ (f_j^i, f_i) \mid i \in I \}$$

and similarly for the columns we define

$$N_J = \{ (f_i^j, f_j) \mid j \in J \}$$

we can say that the profile f_j^i of the rows i is weighted by the marginal probabilities f_i of i .

Definition: The chi-square distance $d^2(i,i')$ between the profiles of two rows i' and i is given by

$$d^2(i,i') = \sum_j \left\{ \frac{1}{f_j} (f_j^i - f_j^{i'})^2 \right\}$$

$$= \sum_j \left\{ \frac{1}{f_j} \left(\frac{f_{ij}}{f_i} - \frac{f_{i'j}}{f_{i'}} \right)^2 \right\} = (1 + \alpha) k(i, j).$$

Furthermore,

$$f_{i_0, j} = \frac{(k(i_0, j))}{K} = \frac{(1 + \alpha) k(i, j)}{K} = \frac{k(i, j)}{K} + \frac{\alpha k(i, j)}{K}$$

i.e.,

$$f_{i_0} = f_{ij} + f_{i'j}$$

Lemma (2):

(preservation of the addition). The chi-square distance between two columns j and j' , is not changed by the addition of the two rows which become one row.

Proof: Let

$$S = \frac{1}{f_i} (f_j^i - f_j^{i'})^2 + \frac{1}{f_{i'}} (f_j^{i'} - f_j^{i''})^2$$

$$= \frac{1}{f_i} \left(\frac{f_{ij}}{f_j} - \frac{f_{i'j}}{f_j} \right)^2 + \frac{1}{f_{i'}} \left(\frac{f_{i'j}}{f_j} - \frac{f_{i''j}}{f_j} \right)^2$$

$$= f_i \left(\frac{f_{ij}}{f_i f_j} - \frac{f_{i'j}}{f_i f_j} \right)^2 + f_{i'} \left(\frac{f_{i'j}}{f_i f_j} - \frac{f_{i''j}}{f_i f_j} \right)^2$$

$$= f_i \left(\frac{f_j^i}{f_j} - \frac{f_j^{i'}}{f_j} \right)^2 + f_{i'} \left(\frac{f_j^{i'}}{f_j} - \frac{f_j^{i''}}{f_j} \right)^2$$

$$= (f_i + f_{i'}) \left(\frac{f_j^{i_0}}{f_j} - \frac{f_j^{i_0'}}{f_j} \right)^2$$

$$= f_{i_0} \left(\frac{f_j^{i_0}}{f_j} - \frac{f_j^{i_0'}}{f_j} \right)^2$$

Since

$$f_j^i = f_j^{i'} = f_j^{i_0} \text{ and } f_j^{i'} = f_j^{i''} = f_j^{i_0'}$$

Similarly the chi-square distance $d^2(j, j')$ between the conditional probabilities f_j^i and $f_j^{i'}$ profiles of the columns j and j' is given by

$$d^2(j, j') = \sum_i \left\{ \frac{1}{f_i} (f_i^j - f_i^{j'})^2 \right\}$$

we can say that the chi-square distance is a weighted Euclidean distance.

This chi-square has the following properties:

Lemma (1):

If the two discrete random variables x and y are both independent of another discrete random variable z then the chi-square distance between their conditional probabilities is zero.

Proof: Consider

$$d^2(x, y) = \sum_z \left\{ \frac{1}{f_z} (f_z^x - f_z^y)^2 \mid z \in Z \right\}$$

$$= \sum_z \left\{ \frac{1}{f_z} \left(\frac{f_{xz}}{f_x} - \frac{f_{yz}}{f_y} \right)^2 \mid z \in Z \right\}$$

$$= \sum_z \left\{ \frac{1}{f_z} \left(\frac{f_x f_z}{f_x} - \frac{f_y f_z}{f_y} \right)^2 \mid z \in Z \right\}$$

$$= \sum_z \left\{ \frac{1}{f_z} (f_z - f_z)^2 \mid z \in Z \right\} = 0$$

If we have a contingency table with two rows i and i' such that $\forall j \in K (i', j) = \alpha k(i, j)$ for some $\alpha \in R^+$ (proportional rows) then it follows that $K(i') = \alpha k(i)$.

Therefore adding i_0 is such that $\forall j \in J$: we have

$$k(i_0, j) = k(i, j) + k(i', j)$$

$$= k(i, j) + \alpha k(i, j)$$

$$= \frac{1}{f_{i_0}} (f_{i_0}^j - f_{i_0}^j)^2$$

Therefore,

$$d^2(j, j') = \sum_i \left\{ \frac{1}{f_i} (f_i^j - f_i^{j'})^2 \mid i \in I \right\}$$

$$= \sum_i \left\{ \frac{1}{f_i} (f_i^j - f_i^{j'})^2 \mid i \in I^* \right\}$$

Note that $\text{Card } I^* = \text{Card } I - 1$

Lemma (3):

Let (X, Y) be the bivariate random variable and A be a vector whose elements are the conditional probability of x with respect to y then $\|A\|^2 \leq \frac{1}{f_X(x)}$ with f_X is the marginal probability of x

Proof:

$$\|A\|^2 = \|f_Y^X\|^2 = \sum_y \left\{ \frac{1}{f_y} (f_y^x)^2 \mid y \in Y \right\}$$

$$\sum_y \left\{ \frac{1}{f_y} \left(\frac{k(x, y)}{k(x)} \right)^2 \mid y \in Y \right\}$$

$$\sum_y \left\{ \frac{1}{f_y} \frac{k(x, y)}{k(x)} \cdot \frac{k(x, y)}{k(x)} \mid y \in Y \right\}$$

$$\sum_y \left\{ \frac{K}{k(y)} \frac{k(x, y)}{k(x)} \cdot \frac{k(x, y)}{k(x)} \mid y \in Y \right\}$$

$$= \frac{K}{k(x)} \sum_y \frac{k(x, y)}{k(y)} \cdot \frac{k(x, y)}{k(x)} \mid y \in Y$$

we know that $\forall y \in Y: \left| \frac{k(x, y)}{k(y)} \right| \leq 1$

$$\|a\|^2 \leq \frac{K}{k(x)} \sum_y \left\{ \frac{k(x, y)}{k(x)} \right\}$$

$$= \frac{K}{k(x)} \frac{k(x)}{k(x)}$$

$$= \frac{K}{k(x)} = \frac{1}{f_X(x)}$$

REFERENCES

- [1] J.P. Benzecri et Collaborateur, *L'analyse des données*, Dunod, 3^e édition. 1980.
- [2] J.P. Benzecri et Collaborateur, *Pratique de l'analyse des données 2*. Abrégé théorique: Etude de cas modèle Dunod, 1980.
- [3] J.P. Nakache, A. Chevalier et V. Maurice, *Exercices Commentés de Mathématiques pour l'analyse statistique des données*, Dunod, Paris, 1981.