

# The Influence Function of the Correlation Indexes in a Two-by-Two Table\*

Giovanni Girone, Fabio Manca, Claudia Marin

University of Bari "Aldo Moro", Bari, Italy

Email: [giovanni.girone@uniba.it](mailto:giovanni.girone@uniba.it), [fabio.manca@uniba.it](mailto:fabio.manca@uniba.it), [claudia.marin@uniba.it](mailto:claudia.marin@uniba.it)

Received 5 October 2014; revised 28 October 2014; accepted 11 November 2014

Copyright © 2014 by authors and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

---

## Abstract

In this paper we examine 5 indexes (the two Yule's indexes, the chi square, the odds ratio and an elementary index) of a two-by-two table, which estimate the correlation coefficient  $\rho$  in a bivariate Bernoulli distribution. We will find the compact expression of the influence functions, which allow the quantification of the effect of an infinitesimal contamination of the probability of any pair of attributes of the bivariate random variable distributed according to the above-mentioned model. We prove that the only unbiased index is the chi square. In order to determine the indexes, which are less sensitive to contamination, we obtain the expressions of three synthetic measures of the influence function, which are the maximum contamination (gross sensitivity error), the mean square deviation and the variance. These results, even if don't allow a definitive assessment of the overall optimum properties of the five indexes, as not all of them are unbiased, nevertheless they allow to appreciating the synthetic entity of the effect of the contaminations in the estimation of the parameter  $\rho$  of the bivariate Bernoulli distribution.

## Keywords

Two-by-Two Table, Influence Function, Correlation Indexes, Gross Sensitivity Error, Mean Square Deviation, Asymptotic Variance

---

## 1. Introduction

In this paper we analyze the influence of a minimal contamination of the bivariate Bernoulli distribution on the values of the index measuring the association in a two-by-two table, having as a scenario the estimation of the correlation parameter of that distribution.

---

\*The paragraphs 1 and 2 are attributed to G. Girone, the paragraphs 3, 4 and 8 to F. Manca and the paragraphs 5, 6 and 7 to C. Marin.

## 2. Bivariate Bernoulli Model

### Selecting a Template

Let us suppose that two dichotomous variables, denoted by  $X$  and  $Y$ , are relevant within a population. These variables take the values 1 and 0, depending on whether one of the dichotomous attributes is present or absent. The corresponding theoretical model is the bivariate Bernoulli distribution [1], reported in **Table 1**.

The mean values of the two variables are

$$E[X] = \gamma + \delta, \tag{1}$$

$$E[Y] = \beta + \delta. \tag{2}$$

The variances of the two variables are

$$Var[X] = (\alpha + \beta)(\gamma + \delta), \tag{3}$$

$$Var[Y] = (\alpha + \gamma)(\beta + \delta). \tag{4}$$

The covariance between the two variables is

$$Cov[X, Y] = \alpha\delta - \beta\gamma. \tag{5}$$

The correlation coefficient is

$$\rho = \frac{\alpha\delta - \beta\gamma}{\sqrt{(\alpha + \beta)(\gamma + \delta)(\alpha + \gamma)(\beta + \delta)}}. \tag{6}$$

## 3. Properties of the Correlation Parameter Estimation

Several indexes, suggested by various authors (Yule, Quetelet and others), are available for the sample estimation of the above-mentioned correlation coefficient. We refer to such indexes as  $R_1, R_2$  etc. For given indexes,  $R_h$ , all variable between  $-1$  and  $+1$ , we must take into account *unbiasedness*, *i.e.*

$$E[R_h] = \rho, \tag{7}$$

*efficiency*, *i.e.*

$$Var[R_h] = \text{minimum}, \tag{8}$$

and the limited influence of limited modification of the model.

With regard to this last fundamental property Hampel [2] in 1974 suggested the *influence function* as a tool for evaluating the effect caused on the value of an index by a minimal contamination of the model. In our case the model is the bivariate Bernoulli distribution, the parameter is the correlation coefficient  $\rho$  and the indexes are those proposed by various authors over time.

Basically the influence function  $If[R]$  referred to the index  $R$  is given by

$$If[R] = \lim_{\varepsilon \rightarrow 0} \frac{R[H_\varepsilon] - R[H]}{\varepsilon}, \tag{9}$$

where  $R[H_\varepsilon]$  is the index computed for the contaminated bivariate Bernoulli distribution  $H_\varepsilon$ ,  $R[H]$  is the

**Table 1.** The bivariate Bernoulli distribution.

Attributes of variable $X$	Attributes of variable $Y$		Total
	0	1	
0	$\alpha$	$\beta$	$\alpha + \beta$
1	$\gamma$	$\delta$	$\gamma + \delta$
Total	$\alpha + \gamma$	$\beta + \delta$	1

index computed for the non-contaminated bivariate Bernoulli distribution and  $\varepsilon$  is the weight of the contamination.

It is easily understood that such a function measures the effect of an infinitesimal contamination of the model on the value of the correlation index [3]. From now on we will denote by  $a$ ,  $b$ ,  $c$  and  $d$  the empirical frequencies of the four cells of the two-by-two table obtained for a sample of  $n$  units.

## 4. Influence Function of the Correlation Indexes

### 4.1. $C$ Index

Let us first consider the elementary index given by

$$C = \frac{a+d-b-c}{a+d+b+c}. \quad (10)$$

A contamination in the cell (0,0) leads to the influence function value

$$If[(0,0), C] = \frac{2(b+c)}{(a+b+c+d)^2}, \quad (11)$$

while a contamination in the cell (1,1) leads to the influence function value

$$If[(1,1), C] = \frac{2(b+c)}{(a+b+c+d)^2}. \quad (12)$$

That is, a contamination in one of the two cells indicating concordance increases the value of the  $C$  index of a quantity, which is proportional to the sum of the frequencies of the two discordance cells.

A contamination in the cell (0,1) leads to the influence function value

$$If[(0,1), C] = -\frac{2(a+d)}{(a+b+c+d)^2}, \quad (13)$$

while a contamination in the cell (1,0) leads to the influence function value

$$If[(1,0), C] = -\frac{2(a+d)}{(a+b+c+d)^2}. \quad (14)$$

That is, a contamination in one of the two cells indicating discordance decreases the value of the  $C$  index of a quantity, which is proportional to the sum of the frequencies of the two concordance cells.

In short, the influence function can be displayed as

$$If[(x,y), C] = \frac{2v}{(a+b+c+d)^2}, \quad (15)$$

in which  $v$ , in the case of a concordance cell, is equal to the sum of the discordance frequencies, while, in the case of a discordance cell, is equal to the sum of the frequencies of discordance cells, changed of sign. In other words, the influence of the contamination for each concordance cell is directly proportional to the sum of the discordant frequencies and vice versa, for each discordance cell, provided that it is positive for the concordance cells and negative for the discordance cells [4].

### 4.2. Yule's $Q$ Index

Let us first consider the 1900 Yule's index [5] given by

$$Q = \frac{ad-bc}{ad+bc}. \quad (16)$$

A contamination in the cell (0,0) leads to the following value of the influence function

$$If[(0,0),Q] = \frac{2bcd}{(ad+bc)^2}, \quad (17)$$

while a contamination in the cell (1,1) leads to the value given by

$$If[(1,1),Q] = \frac{2abc}{(ad+bc)^2}. \quad (18)$$

That is, a contamination in one of the two cells indicating concordance increases the value of the index  $Q$  by a quantity, which is proportional to the product of the frequencies of the three non-contaminated cells.

On the other hand a contamination in the cell (0,1) leads to the value of the influence function given by

$$If[(0,1),Q] = -\frac{2acd}{(ad+bc)^2}, \quad (19)$$

while a contamination in the cell (1,0) leads to the following value of the influence function

$$If[(1,0),Q] = -\frac{2abd}{(ad+bc)^2}. \quad (20)$$

That is, a contamination in one of the two cells indicating discordance decreases the value of the index  $Q$  by a quantity, which is proportional to the product of the frequencies of the three non-contaminated cells.

In short, the influence function can be displayed as

$$If[(x,y),Q] = \frac{2abcd}{v(ad+bc)^2}, \quad (21)$$

in which  $v$  is equal to one of the frequencies with positive sign if it corresponds to  $a$  or to  $d$ , and to one of the frequencies with negative sign if it corresponds to  $b$  or to  $c$ . In other words, the influence of the contamination is inversely proportional to the frequency of the contaminated cell, provided that it is positive for the concordance cells and negative for the discordance cells.

### 4.3. Yule's $Y$ Index

Let us consider now the other index proposed by Yule [6] in 1912,

$$Y = \frac{\sqrt{ad} - \sqrt{bc}}{\sqrt{ad} + \sqrt{bc}}. \quad (22)$$

A contamination in the cell (0,0) leads to the following value of the influence function

$$If[(0,0),Y] = \frac{\sqrt{bcd}}{\sqrt{a}(\sqrt{ad} + \sqrt{bc})^2}, \quad (23)$$

while a contamination in the cell (1,1) leads to the value of the influence function given by

$$If[(1,1),Y] = \frac{\sqrt{abc}}{\sqrt{d}(\sqrt{ad} + \sqrt{bc})^2}. \quad (24)$$

That is, a contamination in one of the two cells indicating concordance increases the value of the index  $Y$  by a quantity proportional to the root of the product of the frequencies of the three non-contaminated cells divided by the root of the frequency of the contaminated cell.

A contamination in the cell (0,1) leads to the value of the influence function given by

$$If[(0,1),Y] = -\frac{\sqrt{acd}}{\sqrt{b}(\sqrt{ad} + \sqrt{bc})^2}, \quad (25)$$

while a contamination in the cell (1,0) leads to the following value of the influence function

$$If[(0,1), Y] = -\frac{\sqrt{abd}}{\sqrt{c}(\sqrt{ad} + \sqrt{bc})^2}. \quad (26)$$

That is, a contamination in one of the two cells indicating discordance decreases the value of the index  $Y$  by a quantity proportional to the root of the product of the frequencies of the three non-contaminated cells divided by the root of the frequency of the contaminated cell.

In short, the influence function can be displayed as

$$If[(x, y), Y] = \frac{\sqrt{abcd}}{v(\sqrt{ad} + \sqrt{bc})^2}, \quad (27)$$

in which  $v$  is equal to one of the frequencies with positive sign if  $a$  or  $d$  and to one of the frequencies with negative sign if  $b$  or  $c$ . In other words, the influence of the contamination is inversely proportional to the frequency of the contaminated cell, provided that it is positive for the concordance cells and negative for the discordance cells.

#### 4.4. The Chi Square Index

Let us examine the chi square index

$$\chi^2 = \frac{ad - bc}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}. \quad (28)$$

A contamination in the cell (0,0) leads to the influence function value

$$If[(0,0), \chi^2] = \frac{(b+d)(c+d)[(b+c)(ad+bc) + 2bc(a+d)]}{2[(a+b)(a+c)(b+d)(c+d)]^{3/2}}, \quad (29)$$

while a contamination in the cell (1,1) leads to the influence function value

$$If[(1,1), \chi^2] = \frac{(a+b)(a+c)[(b+c)(ad+bc) + 2bc(a+d)]}{2[(a+b)(a+c)(b+d)(c+d)]^{3/2}}. \quad (30)$$

That is, a contamination in one of the two cells indicating concordance increases the value of the  $\chi^2$  by a quantity which is proportional to the product of the sums of the frequency of the other concordance cell with each of the frequencies of the discordance cell.

A contamination in the cell (0,1) leads to the influence function value

$$If[(0,1), \chi^2] = -\frac{(a+c)(c+d)[(a+d)(ad+bc) + 2ad(b+c)]}{2[(a+b)(a+c)(b+d)(c+d)]^{3/2}}, \quad (31)$$

while a contamination in the cell (1,0) leads to the influence function value

$$If[(1,0), \chi^2] = -\frac{(a+b)(b+d)[(a+d)(ad+bc) + 2ad(b+c)]}{2[(a+b)(a+c)(b+d)(c+d)]^{3/2}}. \quad (32)$$

That is, a contamination in one of the two cells indicating discordance decreases the value of the  $\chi^2$  by a quantity which is proportional to the product of the sums of the frequency of the other discordance cell with each of the frequencies of the concordance cells.

It is impossible to have a unique expression of the influence function as we had for the other indexes, because the expressions for the contaminated concordance cells differ from those related to the discordance ones.

#### 4.5. Odds Ratio, $\theta$

Let us now examine the odds ratio index

$$\theta = \frac{ad}{bc}. \quad (33)$$

A contamination in the cell (0,0) leads to the influence function value

$$If[(0,0),\theta] = \frac{d}{bc}, \quad (34)$$

while a contamination in the cell (1,1) leads to the influence function value

$$If[(1,1),\theta] = \frac{a}{bc}. \quad (35)$$

That is, a contamination in one of the two cells indicating concordance increases the value of the index  $\theta$  by a quantity that is proportional to the frequency of the other concordance cell and inversely proportional to the product of the frequencies of the discordance cells.

A contamination in the cell (0,1) leads to the influence function value

$$If[(0,1),\theta] = -\frac{ad}{b^2c}, \quad (36)$$

while a contamination in cell (1,0) leads to the influence function value

$$If[(1,0),\theta] = -\frac{ad}{bc^2}. \quad (37)$$

That is, a contamination in one of the two cells indicating discordance decreases the value of the index  $\theta$  by a quantity, which is proportional to the product of the frequencies of the concordance cell and inversely proportional to the product of the square of the frequency of the contaminated cell multiplied by the frequency of the other discordance cell.

In short, the influence function can be displayed as

$$If[(x,y),\theta] = \frac{ad}{bcv} = \frac{\theta}{v}, \quad (38)$$

in which  $v$  is equal to the frequency of the contaminated cell, provided that the sign is positive in case of contamination in a concordance cell, and negative in case of contamination in a discordance cell.

## 5. Unbiasedness of the Indexes

It must be reminded that an index  $R_h$  is unbiased if

$$E(R_h) = \rho, \quad (39)$$

let us examine now the unbiasedness of every index.

### 5.1. The Chi Square Index

The index

$$\chi^2 = \frac{ad - bc}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}, \quad (40)$$

has the mean

$$E(\chi^2) = E\left(\frac{ad - bc}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}\right) = \rho, \quad (41)$$

and therefore it is unbiased.

### 5.2. Other Indexes

Indexes  $Q$ ,  $Y$ ,  $\theta$  e  $C$  are biased.

It has to be said that 3 of the considered indexes ( $Q$ ,  $Y$  and  $\theta$ ) are functionally related, as it is shown below:

$$Q = \frac{2Y}{1+Y^2}, \quad Y = \frac{1+\sqrt{1-Q^2}}{Q} \quad \text{for } -1 < Q < 0 \quad \text{and} \quad Y = \frac{1-\sqrt{1-Q^2}}{Q} \quad \text{for } 0 < Q < 1, \quad Q = \frac{\theta-1}{\theta+1},$$

$$\theta = \frac{1+Q}{1-Q}, \quad \theta = \frac{(1+G)^2}{(1-G)^2}, \quad G = \frac{1-\sqrt{\theta}}{1+\sqrt{\theta}}, \quad \text{for } 0 < \theta < 1 \quad \text{and} \quad G = \frac{1+\sqrt{\theta}}{1-\sqrt{\theta}}, \quad \text{for } 1 < \theta < \infty.$$

The other 2 indexes ( $C$  and  $\chi^2$ ) are not functionally explainable with themselves nor with the above-mentioned ones. The 5 indexes estimate functions of the  $\rho$  parameter. More exactly 4 of these indexes ( $C$ ,  $Q$ ,  $Y$  and  $\chi^2$ ) are estimators of increasing functions of this parameter and, in particular in the points  $-1$ ,  $0$  and  $+1$ , these functions coincide with the argument. So the index  $\theta$  can easily lead to the 2 Yule's indexes achieving again its characteristics.

## 6. Influences of the Indexes

Since that the effects of contaminations in the various cells are balanced, it is necessary to evaluate their overall influence regardless of the sign. This can be done considering the *maximum of the absolute values* of the influence or the *mean absolute deviation* or the *variance* of the said values [7].

### 6.1. Maximum of the Absolute Values of the Influence Function (Gross Sensitivity Error)

#### 6.1.1. C Index

As, regardless of the sign, the influence function is equal to

$$If(v, C) = \frac{a+d}{(a+b+c+d)^2} \quad \text{or} \quad If(v, C) = \frac{b+c}{(a+b+c+d)^2}, \quad (42)$$

the maximum of the influence function is therefore

$$GSE(C) = \frac{2(a+d)(b+c)}{\min(a+d, b+c)(a+b+c+d)^2}. \quad (43)$$

#### 6.1.2. Yule's Q Index

As, regardless of the sign, the influence function is equal to

$$If(v, Q) = \frac{2abcd}{v(ad+bc)^2}, \quad (44)$$

in which  $v$  is one of the four frequencies of the table, the maximum of the influence function is obtained for  $\min(v)$ ; it is therefore

$$GSE(Q) = \frac{2abcd}{\min(a, b, c, d)(ad+bc)^2}. \quad (45)$$

#### 6.1.3. Yule's Y Index

As, regardless of the sign, the influence function is equal to

$$If(v, Y) = \frac{\sqrt{abcd}}{v(\sqrt{ad} + \sqrt{bc})^2}, \quad (46)$$

in which  $v$  is one of the four frequencies of the table, the maximum of the influence function is obtained for  $\min(v)$ ; it is therefore

$$GSE(Y) = \frac{\sqrt{abcd}}{\min(a, b, c, d)(\sqrt{ad} + \sqrt{bc})^2}. \quad (47)$$

### 6.1.4. Chi Square Index

An empirical analysis allows to assess that the maximum absolute value of the influence function is obtained in correspondence of the minimum frequency. Thus,

$$GSE(\chi^2) = If[(0,0), \chi^2] \text{ if } a < \min(b, c, d), \quad (48)$$

$$GSE(\chi^2) = If[(1,1), \chi^2] \text{ if } d < \min(a, b, c), \quad (49)$$

$$GSE(\chi^2) = If[(0,1), \chi^2] \text{ if } b < \min(a, c, d), \quad (50)$$

$$GSE(\chi^2) = If[(1,0), \chi^2] \text{ if } c < \min(a, b, d). \quad (51)$$

### 6.1.5. Odds Ratio

As, regardless of the sign, the influence function is equal to

$$If(v, \theta) = \frac{ad}{vbc}, \quad (52)$$

in which  $v$  is one of the four frequencies of the table, the maximum of the influence function is obtained for  $\min(v)$ ; it is therefore

$$GSE(\theta) = \frac{ad}{\min(a, b, c, d)bc}. \quad (53)$$

## 6.2. Variability of Influence Functions: Mean Absolute Deviation

### 6.2.1. C Index

A few algebraic steps allow us to obtain

$$MD(C) = \frac{4(a+d)(b+c)}{(a+b+c+d)^3}. \quad (54)$$

### 6.2.2. Yule's Q Index

$$MD(Q) = \frac{8adbc}{(ad+bc)^2}. \quad (55)$$

### 6.2.3. Yule's Y Index

$$MD(Y) = \frac{4\sqrt{adbc}}{ad+bc+2\sqrt{adbc}}. \quad (56)$$

### 6.2.4. Chi Square

$$MD(\chi^2) = \frac{[2bc(a+d) + (b+c)(ad+bc)][2ad(b+c) + (a+d)(ad+bc)]}{[(a+b)(a+c)(b+d)(c+d)]^{3/2}}. \quad (57)$$

### 6.2.5. Odds Ratio

$$MD(\theta) = \frac{4ad}{bc}. \quad (58)$$

It can be seen that the mean deviation for all indexes is a symmetric function either of the concordant fre-



quencies or of the discordant frequencies.

### 6.3. Variability of the Influence Function Asymptotic Variance (A.S.V.)

#### 6.3.1. *C* Index

Let us consider the asymptotic variance of the indexes. A few algebraic steps lead us to the following expression

$$ASV(C) = \frac{4(a+d)(b+c)}{(a+b+c+d)^4}. \quad (59)$$

#### 6.3.2. Yule's *Q* Index

$$ASV(Q) = \frac{4adbc[ad(b+c)+bc(a+d)]}{(ad+bc)^4}. \quad (60)$$

#### 6.3.3. Yule's *Y* Index

$$ASV(Y) = \frac{ad(b+c)+bc(a+d)}{(ad+bc+2\sqrt{adbc})^2}. \quad (61)$$

#### 6.3.4. Chi Square Index

$$ASV(\chi^2) = \frac{\{(a+d)(b+c)(6adbc+b^2c^2+a^2d^2)+4adbc[(a+d)^2+(b+c)^2]\}}{[2(a+b)(a+c)(b+d)(c+d)]^2}. \quad (62)$$

#### 6.3.5. Odds Ratio

$$ASV(\theta) = \frac{ad[ad(b+c)+bc(a+d)]}{(bc)^3}. \quad (63)$$

It can be seen that the asymptotic variance is a symmetric function of the concordance and discordance frequencies as well.

## 7. Example

Let us consider a practical example in which 1071 persons are classified on 2 dichotomic characters: "does he/she smoke" and "is he/she suffering from bronchitis?" both with yes or no response (see [Table 2](#)).

There were 1071 cases of which 135 smoke and have bronchitis and 547 don't smoke and don't have bronchitis.

Elementary indexes	<i>C</i>	<i>Q</i>	<i>Y</i>	$\chi^2$	$\theta$
G.S.E.	0.00119	0.00399	0.00232	0.00179	0.02473
M.D.	0.00086	0.00152	0.00089	0.00076	0.00942
A.S.V.	0.0000008	0.0000035	0.0000012	0.0000009	0.0001338
S.D.	0.00090	0.00186	0.00109	0.00094	0.01157

As it can be noticed, between the 4 indexes whose values go between  $-1$  and  $+1$ , the ones which are less sensitive to contamination are *C* and chi square indexes; on the other hand, the more sensitive ones are Yule's indexes, *Q* and *Y*. The greater sensitivity of the odds ratio is due to the fact that such index measures a function of the correlation of the model that goes in the range from 1 to  $\infty$ .

**Table 2.** Smoke versus bronchitis.

Smoke	Bronchitis		Total
	Yes	No	
Yes	135	287	422
No	102	547	649
Total	237	834	1071

Source: Survey at the University Hospital of Bari, Department of Pulmonology.

## 8. Conclusions

In this paper we analyzed the indexes of a two-by-two table, which allow the estimation of the correlation coefficient  $\rho$  in the bivariate Bernoulli model. More precisely, we considered the two Yule's indexes, the chi square, the odds ratio and a further elementary index. We obtained, for these indexes, the compact expressions of the influence functions, which allow the quantification of the effect of an infinitesimal contamination of the probability of any pair of attributes of the bivariate random variable distributed according to the above-mentioned model.

In order to determine the indexes which are less sensitive to contamination, we obtained the expressions of three synthetic measures of the influence function, specifically the maximum contamination (gross sensitivity error), the mean absolute deviation and the variance. These expressions, even if don't allow a definitive assessment of the overall optimum properties of the five indexes considered, as not all of them are unbiased, nevertheless they allow to appreciating the synthetic entity of the effect of the contaminations in the estimation of the parameter  $\rho$  of the bivariate Bernoulli model.

## References

- [1] Barnard, G.A. (1981) Two by Two ( $2 \times 2$ ) Tables. *Encyclopedia of Statistical Sciences*, **9**, 367-372.
- [2] Hampel, F.R. (1974) The Influence Curve and Its Role in Robust Estimation. *Journal of the American Statistical Association*, **69**, 383-393. <http://dx.doi.org/10.1080/01621459.1974.10482962>
- [3] Kendall, M.G. and Stuart, A. (1977) *The Advanced Theory of Statistics*. Vol. 2, C. Griffin, London, 566-571.
- [4] Pearson, K. (1904) On the Theory of Contingency and Its Relation to Association and Normal Correlation. *Biometric Series*, Drapers' Co. Memoirs, London.
- [5] Yule, G.U. (1900) On the Association of Attributes in Statistics. *Philosophical Transaction*, **194**, 257. <http://dx.doi.org/10.1098/rsta.1900.0019>
- [6] Yule, G.U. (1912) On the Methods of Measuring Association between Two Attributes. *Journal of the Royal Statistical Society*, **75**, 579. <http://dx.doi.org/10.2307/2340126>
- [7] Yule, G.U. and Kendal, M.G. (1958) *An Introduction to the Theory of Statistics*. C. Griffin, London, 271-272.