

## On characterizing classical bivariate normality via regression functions

Barry C. Arnold <sup>a</sup> and B.G. Manjunath <sup>b</sup>

<sup>a</sup> Department of Statistics, University of California, Riverside, USA. <sup>b</sup> School of Mathematics and Statistics, University of Hyderabad, Hyderabad

### ARTICLE HISTORY

Compiled December 12, 2021

Received 09 July 2021; Accepted 20 September 2021

### ABSTRACT

The possibility of characterizing classical bivariate normality via regression conditions is discussed. Some additional assumptions are required. It has been suggested that a further assumption of constant conditional variance would suffice, but even this is not adequate. It is shown that an additional assumption of conditional normality of  $X$  given  $Y = y$  for all  $y$  will yield the desired characterization. Related characterization problems are also considered. The analogous trivariate problem is discussed but is unresolved.

### KEYWORDS

Bivariate normal; Conditional normal; Linear regression; Moment generating function; Variance-covariance matrix

## 1. Introduction

Consider bivariate densities such that for every  $y$

$$\int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx = a + by$$

and for every  $x$

$$\int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy = c + dx.$$

Regression functions of this type are familiar since they are encountered in the case of a classical bivariate normal distribution. It might be tempting to speculate that the existence of linear regression functions might characterize the classical model.

To investigate this possibility, we may recall a construction suggested in Wesolowski [7] (see also Arnold and Wesolowski [2]). Begin with a joint density  $f(x, y)$  corresponding to a particular classical bivariate normal distribution. Let  $f_1$  and  $f_2$  be two distinct bounded densities on the interval  $(-1, 1)$  each having mean 0 and variance  $1/2$ . Let  $c$  be a small positive number and consider the function

$$f^*(x, y) = f(x, y) + c[f_1(x) - f_2(x)][f_1(y) - f_2(y)]. \quad (1)$$

Provided that  $c$  is chosen to be small enough to ensure that  $f^*$  is always positive, then  $f^*$  is a non-normal bivariate density, but it is easy to verify that it has the same marginals and the same regression functions and conditional variances as does the density  $f$ . For example: one can consider the densities

**Example 1.1.**

$$\begin{aligned} f_1(x) &= \frac{3}{8} - \frac{|x|}{4}; \quad -1 < x < 1, \\ f_2(x) &= \frac{39}{48} - \frac{15|x|^2}{16}; \quad -1 < x < 1. \end{aligned}$$

Spanos [5], using a result of Nimmo-Smith [6], argued that linear regressions and the imposition of the additional condition that the conditional variance of  $X$  given  $Y = y$  is constant; is enough to guarantee that we have a classical bivariate normal model. Examples of the form (1) contradict this claim.

Some additional condition is required. Bhattacharyya [3] suggested several conditions involving moment assumptions and assumptions of normality of both sets of conditional distributions to ensure that the classical bivariate model is in place. See also Castillo and Galambos [4] for other conditions. A good summary of characterization results involving conditional distributions may be found in Arnold, Castillo and Sarabia [1].

In Section 2, we will add one additional distributional assumption to the list used by Spanos and we obtain a characterization of classical bivariate normality.

**Note 1.** If we merely assume linear regressions and finite variances for  $X$  and  $Y$ , then with no other distributional assumptions (and some tedious algebra) we can verify that the variance-covariance matrix of  $(X, Y)$  is of the form

$$\Sigma = \begin{pmatrix} \sigma_X^2 & d\sigma_X^2 \\ d\sigma_X^2 & \frac{d}{b}\sigma_X^2 \end{pmatrix}. \quad (2)$$

## 2. The characterization

One might consider adding an assumption of normality of one or perhaps both marginals to Spanos' list of conditions to hopefully guarantee that the joint density is bivariate normal. Unfortunately this will not work since our example in (1) has normal marginals in addition to satisfying Spanos' conditions.

We will instead augment Spanos' conditions in a different fashion and ask whether the following two conditions are sufficient to imply that  $(X, Y)$  has a classical bivariate normal distribution.

$$X|Y = y \sim \text{Normal}(a + by, k), \quad \forall y \quad (3)$$

and

$$E(Y|X = x) = c + dx, \quad \forall x. \quad (4)$$

These conditions are clearly stronger than the Spanos conditions. In the proof we make use of the following observations. A random variable  $X$  has a normal distribution if  $f_X(x) \propto e^{c_1x + c_2x^2}$  for some real  $c_1$  and negative  $c_2$ . Analogously  $X$  is normally distributed if its moment generating function is of the form  $M_X(t) = e^{d_1t + d_2t^2}$  for some real  $d_1$  and positive  $d_2$ .

**Lemma 2.1.** *If conditions (3) and (4) hold then necessarily  $Y$  has a normal distribution and consequently,  $(X, Y)$  has a classical bivariate normal distribution.*

**Proof.** The density of  $Y$  to be denoted by  $g(y)$  will be obtained by solving the following system of equations

$$\int_{-\infty}^{\infty} [c + dx - y] \frac{e^{-(x-a-by)^2/2k}}{\sqrt{k2\pi}} g(y) dy = 0, \quad \forall x \in (-\infty, \infty).$$

It is clear that if we can solve the problem for the case in which  $k = 1$ , then the general solution can be readily obtained.

With  $k = 1$  our equation is of the form:

$$\int_{-\infty}^{\infty} [c + dx - y] \frac{e^{-(x-a-by)^2/2}}{\sqrt{2\pi}} g(y) dy = 0, \quad \forall x \in (-\infty, \infty)$$

or

$$\int_{-\infty}^{\infty} [c + dx - y] e^{-(x-a-by)^2/2} g(y) dy = 0, \quad \forall x \in (-\infty, \infty).$$

Equivalently

$$\int_{-\infty}^{\infty} [c + dx - y] e^{bxy} \left\{ e^{-b^2y^2/2} e^{-aby} g(y) \right\} dy = 0, \quad \forall x \in (-\infty, \infty).$$

or

$$\int_{-\infty}^{\infty} [c + dx - y]e^{bxy}h(y)dy = 0, \quad \forall x \in (-\infty, \infty), \quad (5)$$

where

$$h(y) = \frac{\{e^{-b^2y^2/2}e^{-aby}g(y)\}}{\int_{-\infty}^{\infty} \{e^{-b^2y^2/2}e^{-aby}g(y)\} dy}.$$

Replace  $bx$  by  $t$  in (5) to get

$$\int_{-\infty}^{\infty} [c + d_1t - y]e^{ty}h(y)dy = 0, \quad \forall t \in (-\infty, \infty), \quad (6)$$

where  $d_1 = d/b$ .

Denote the m.g.f. of  $h$  by  $M(t)$ .

Thus  $M(t)$  satisfies the following differential equation

$$[c + d_1t]M(t) = M'(t)$$

with solution  $M(t) = \exp[ct + d_2t^2]$ . Thus  $h$  is a normal density from

which it follows that  $g$  is a normal density. □

Conditions (3) and (4) involve 5 parameters  $a, b, c, d$  and  $k$ . Having decided that  $(X, Y)$  has a classical bivariate normal distribution, it is not difficult to identify the means, variances and covariance of its distribution as functions of  $a, b, c, d$  and  $k$ . Thus we have

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N \left( \begin{pmatrix} \frac{a+bc}{1-bd} \\ \frac{c+ad}{1-bd} \end{pmatrix}, \begin{pmatrix} \frac{k}{1-bd} & \frac{dk}{b(1-bd)} \\ \frac{dk}{1-bd} & \frac{d}{k(1-bd)} \end{pmatrix} \right)$$

Note that for this to be a valid distribution we must have  $0 < bd < 1$ .

### 3. A related characterization

Consider the following conditions

$$X|Y = y \sim \text{Normal}(a + by, k) \quad \forall y \quad (7)$$

and

$$E(Y|X = x) = \delta(x) \quad \forall x. \quad (8)$$

**Lemma 3.1.** *If conditions (7) and (8) hold the distribution of  $Y$  is uniquely determined*

*and so is the joint distribution of  $(X, Y)$ .*

**Proof.** The density of  $Y$ , to be denoted by  $g(y)$ , will be obtained by solving

the following system of equations

$$\int_{-\infty}^{\infty} [\delta(x) - y] \frac{e^{-(x-a-by)^2/2k}}{\sqrt{k}2\pi} g(y) dy = 0, \quad \forall x \in (-\infty, \infty).$$

It is clear that if we can solve the problem for the case in which  $k = 1$ , then the general solution can be readily obtained.

With  $k = 1$  our equation is of the form:

$$\int_{-\infty}^{\infty} [\delta(x) - y] \frac{e^{-(x-a-by)^2/2}}{\sqrt{2\pi}} g(y) dy = 0, \quad \forall x \in (-\infty, \infty).$$

or

$$\int_{-\infty}^{\infty} [\delta(x) - y] e^{-(x-a-by)^2/2} g(y) dy = 0, \quad \forall x \in (-\infty, \infty).$$

Equivalently

$$\int_{-\infty}^{\infty} [\delta(x) - y] e^{bxy} \left\{ e^{-b^2y^2/2} e^{-aby} g(y) \right\} dy = 0, \quad \forall x \in (-\infty, \infty).$$

or

$$\int_{-\infty}^{\infty} [\delta(x) - y] e^{bxy} h(y) dy = 0, \quad \forall x \in (-\infty, \infty), \quad (9)$$

where

$$h(y) = \frac{\{e^{-b^2y^2/2} e^{-aby} g(y)\}}{\int_{-\infty}^{\infty} \{e^{-b^2y^2/2} e^{-aby} g(y)\} dy}.$$

Replace  $bx$  by  $t$  in (9) to get

$$\int_{-\infty}^{\infty} [\delta(t/b) - y] e^{ty} h(y) dy = 0, \quad \forall t \in (-\infty, \infty),$$

where  $d_1 = d/b$ .

Denote the m.g.f. of  $h$  by  $M(t)$ .

Thus  $M(t)$  satisfies the following differential equation

$$[\delta(t/b)]M(t) = M'(t)$$

with solution  $M(t) = \exp[\int_0^t \delta(s/b) ds]$ . Thus  $h$  is determined and thus  $g(y)$  is determined. □

#### 4. In three dimensions

Suppose that the random variable  $(X, Y, Z)$  has linear regression functions, i.e.,

$$E(X|Y = y, Z = z) = \alpha_1 + \gamma_1 y + \delta_1 z, \quad (10)$$

$$E(Y|X = x, Z = z) = \alpha_2 + \beta_1 x + \delta_2 z, \quad (11)$$

$$E(Z|X = x, Y = y) = \alpha_3 + \beta_2 x + \gamma_2 y, \quad (12)$$

By considering the case in which the means are zero and second moments exist, and denoting the variance-covariance matrix by  $\Sigma$  with elements  $\{\sigma_{ij}\}_{i=1,j=1}^{3,3}$  and setting  $\sigma_{11} = 1$ , we are able to completely identify  $\Sigma$  using (10)-(12) by solving the following system of 5 equations.

$$\begin{pmatrix} -1 & 0 & \gamma_1 & \delta_1 & 0 \\ -1 & 0 & \delta_2 & 0 & 0 \\ \gamma_2 & -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & \gamma_1 & \delta_1 \\ \beta_2 & 0 & \gamma_2 & -1 & 0 \end{pmatrix} \begin{pmatrix} \sigma_{12} \\ \sigma_{13} \\ \sigma_{22} \\ \sigma_{23} \\ \sigma_{33} \end{pmatrix} = \begin{pmatrix} 0 \\ -\beta_1 \\ -\beta_2 \\ 0 \\ 0 \end{pmatrix}. \quad (13)$$

It follows that

$$\begin{pmatrix} \sigma_{12} \\ \sigma_{13} \\ \sigma_{22} \\ \sigma_{23} \\ \sigma_{33} \end{pmatrix} = \begin{pmatrix} -1 & 0 & \gamma_1 & \delta_1 & 0 \\ -1 & 0 & \delta_2 & 0 & 0 \\ \gamma_2 & -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & \gamma_1 & \delta_1 \\ \beta_2 & 0 & \gamma_2 & -1 & 0 \end{pmatrix}^{-1} \begin{pmatrix} 0 \\ -\beta_1 \\ -\beta_2 \\ 0 \\ 0 \end{pmatrix}. \tag{14}$$

$$= \begin{pmatrix} \frac{\delta_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & -\frac{\gamma_1 + \delta_1\gamma_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & 0 & 0 & \frac{\delta_1\delta_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} \\ \frac{\delta_2\gamma_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & -\frac{\gamma_2(\gamma_1 + \delta_1\gamma_2)}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & -1 & 0 & \frac{\delta_1\delta_2\gamma_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} \\ \frac{1}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & \frac{\delta_1\beta_1 - 1}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & 0 & 0 & \frac{\delta_1}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} \\ \frac{\delta_2\beta_2 + \gamma_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & -\frac{\gamma_1\beta_2 + \gamma_2}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} & 0 & 0 & \frac{\delta_2 - \gamma_1}{\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2} \\ \frac{\delta_2\gamma_2 - \gamma_1(\delta_2\beta_2 + \gamma_2)}{\delta_1(\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2)} & \frac{\gamma_1^2\beta_2 - \delta_1\gamma_2^2}{\delta_1(\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2)} & -\frac{1}{\delta_1} & \frac{1}{\delta_1} & \frac{\gamma_1^2 - \gamma_1\delta_2 + \delta_1\delta_2\gamma_2}{\delta_1(\gamma_1 + \delta_2(\delta_1\beta_2 - 1) + \delta_1\gamma_2)} \end{pmatrix} \begin{pmatrix} 0 \\ -\beta_1 \\ -\beta_2 \\ 0 \\ 0 \end{pmatrix}. \tag{15}$$

However, we know that having linear regressions as in (10),(11) and (12) is not enough to imply that  $(X, Y, Z)$  has a classical trivariate normal distribution, since a three dimensional version of the density in (1) is readily constructed. Paralleling the discussion of the two-dimensional case, we might consider imposing the additional condition that

$$X|Y = y, Z = z \sim N(\alpha_1 + \gamma_1 y + \delta_1 z, k), \quad \forall y, z \in (-\infty, \infty). \tag{16}$$

This assumption will lead to differential equations involving the joint mgf of  $(Y, Z)$  but we have been unable to identify this mgf and, consequently, we cannot decide whether the joint density of  $(X, Y, Z)$  is of the classical trivariate normal form.

Of course, if we made the following three assumptions (which include linear regression conditions),

$$X|Y = y, Z = z \sim N(\alpha_1 + \gamma_1 y + \delta_1 z, k_1), \quad \forall y, z \in (-\infty, \infty), \tag{17}$$

$$Y|X = x, Z = z \sim N(\alpha_2 + \beta_1 x + \delta_2 z, k_2), \quad \forall x, z \in (-\infty, \infty), \tag{18}$$

$$Z|X = x, Y = y \sim N(\alpha_3 + \beta_2 x + \gamma_2 y, k_3), \quad \forall x, y \in (-\infty, \infty), \tag{19}$$

then we would have a three dimensional normal conditionals model and it can be argued that consequently  $(X, Y, Z)$  must have a classical normal distribution.

Since, in two dimensions we only needed to add one conditional normality assumption to our linear regression assumptions, it seems unlikely that in three dimensions we would need all three conditional normality assumptions as in (17),(18) and (19). The problem is currently unresolved.

**Acknowledgement(s)**

The second author’s research was sponsored by the Institution of Eminence (IoE), University of Hyderabad (UoH-IoE-RC2-21-013).

**References**

- [1] Arnold, B.C., Castillo, E., and Sarabia, J.M., Conditional Specification of Statistical Models, Springer Series in Statistics, New York (1999).
- [2] Arnold, B.C. and Wesolowski, J., Multivariate distributions with Gaussian conditional structure, in book: Stochastic Processes and Functional Analysis: in celebration of M.M. Rao's 65th birthday, 45–59, (1997).
- [3] Bhattacharyya, A., On some sets of sufficient conditions leading to the normal bivariate distribution, Sankhya, 6, 399–406, (1943).
- [4] Castillo, E. and Galambos, J., Conditional distributions and the bivariate normal distributions, Metrika, 36, 209–214, (1989).
- [5] Spanos, A., On normality and the linear regression model, Econometric Reviews, 14(2), 195–203, (1995).
- [6] Nimmo-Smith, I., Linear regressions and sphericity, Biometrika, 66, 390–392, (1979).
- [7] Wesolowski, J., Gaussian conditional structure of the second order and the Kagan classification of multivariate distributions, Journal of Multivariate Analysis, 39, 79–86, (1991).