



On an independence test approach to the goodness-of-fit problem



Ludwig Baringhaus*, Daniel Gaigall

Institut für Mathematische Stochastik, Leibniz Universität Hannover, Postfach 60 09, D-30060 Hannover, Germany

ARTICLE INFO

Article history:

Received 10 October 2014

Available online 27 May 2015

AMS 2000 subject classifications:

62G10

62G09

Keywords:

Goodness-of-fit test

Independence test

Parametric bootstrap

Vapnik–Červonenkis class

U-process

Gamma distribution

Inverse Gaussian distribution

ABSTRACT

Let X_1, \dots, X_n be independent and identically distributed random variables with distribution \mathbb{F} . Assuming that there are measurable functions $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ characterizing a family \mathcal{F} of distributions on the Borel sets of \mathbb{R} in the way that the random variables $f(X_1, X_2), g(X_1, X_2)$ are independent, if and only if $\mathbb{F} \in \mathcal{F}$, we propose to treat the testing problem $H: \mathbb{F} \in \mathcal{F}, K: \mathbb{F} \notin \mathcal{F}$ by applying a consistent nonparametric independence test to the bivariate sample variables $(f(X_i, X_j), g(X_i, X_j)), 1 \leq i, j \leq n, i \neq j$. A parametric bootstrap procedure needed to get critical values is shown to work. The consistency of the test is discussed. The power performance of the procedure is compared with that of the classical tests of Kolmogorov–Smirnov and Cramér–von Mises in the special cases where \mathcal{F} is the family of gamma distributions or the family of inverse Gaussian distributions.

© 2015 Elsevier Inc. All rights reserved.

1. Introduction

Let \mathcal{G} be some general nonparametric class of non-degenerate distributions on the Borel sets of \mathbb{R} , and let $\emptyset \neq \mathcal{F} = \{\mathbb{F}(\cdot; \vartheta); \vartheta \in \Theta\} \subsetneq \mathcal{G}$ be a parametric subfamily indexed by some parameter $\vartheta \in \Theta$, where $\Theta \neq \emptyset$ is a subset of \mathbb{R}^d , say. Let X_1, \dots, X_n, \dots be real valued, independent and identically distributed random variables with unknown distribution $\mathbb{F} \in \mathcal{G}$. Let us express the fact that X_1 has distribution (function) \mathbb{F} by $X_1 \sim \mathbb{F}$. On the basis of X_1, \dots, X_n we consider testing the hypothesis

$$H: \mathbb{F} \in \mathcal{F} \tag{1.1}$$

against the general alternative $K: \mathbb{F} \in \mathcal{G} \setminus \mathcal{F}$. Let us assume that there are measurable functions $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ characterizing the family \mathcal{F} in the way that the random variables $f(X_1, X_2), g(X_1, X_2)$ are independent, if and only if $\mathbb{F} \in \mathcal{F}$.

Example 1.1. (1) The independence of $X_1 - X_2$ and $X_1 + X_2$ characterizes the family of normal distributions [5].

(2) Let the X_i be non-negative. Then $\frac{X_1}{X_1 + X_2}$ and $X_1 + X_2$ are independent, if and only if \mathcal{F} is the family of Gamma distributions $G(\alpha, \lambda)$ with shape parameter $\alpha > 0$ and scale parameter $\lambda > 0$ [22].

(3) Let the X_i be positive with finite moments $E(X_i^2)$ and $E(X_i^{-1})$. Then $\bar{X} = (X_1 + X_2)/2$ and $V = \frac{1}{2} \left(\frac{1}{X_1} + \frac{1}{X_2} \right) - \frac{1}{\bar{X}}$ are independent if and only if \mathcal{F} is the family of inverse Gaussian distributions $IG(\mu, \lambda)$ with densities given by $\sqrt{\frac{\lambda}{2\pi}} x^{-3/2} \exp\left(-\frac{\lambda(x-\mu)}{2\mu^2 x}\right), x > 0$, for parameter values $\mu > 0$ and $\lambda > 0$ [19].

* Corresponding author.

E-mail addresses: lbaring@stochastik.uni-hannover.de (L. Baringhaus), gaigall@stochastik.uni-hannover.de (D. Gaigall).

- (4) Let $0 < X_1 < 1$. Then the independence of $\frac{1-X_1}{1-X_1X_2}$ and $1 - X_1X_2$ characterizes the family of beta distributions [30].
- (5) Let $X_1 > 0$ and let the distribution function of X_1 be strictly increasing. Then $\min(X_1, X_2)$ and $|X_1 - X_2|$ are independent if and only if X_1 has an exponential distribution [12, Theorem 3.3.1].

Given independent and identically distributed bivariate random vectors $(Y_1, Z_1), \dots, (Y_n, Z_n)$ with absolutely continuous distribution, the hypothesis of independence of Y_j and Z_j can simply and consistently be tested by using the Hoeffding–Blum–Kiefer–Rosenblatt independence criterion, that is by rejecting the hypothesis of independence for large values of

$$T_n = n \int (H_n(y, z) - F_n(y)G_n(z))^2 dH_n(y, z),$$

where

$$H_n(y, z) = \frac{1}{n} \sum_{j=1}^n I(Y_j \leq y, Z_j \leq z), \quad (y, z) \in \mathbb{R}^2,$$

$$F_n(y) = H_n(y, \infty) = \frac{1}{n} \sum_{j=1}^n I(Y_j \leq y), \quad y \in \mathbb{R},$$

$$G_n(z) = H_n(\infty, z) = \frac{1}{n} \sum_{j=1}^n I(Z_j \leq z), \quad z \in \mathbb{R},$$

are the empirical distribution functions of the joint and the marginal sample variables. Defining for each $1 \leq j \leq n$

$$N_1(j) = \sum_{v=1}^n I(Y_v \leq Y_j, Z_v \leq Z_j), \quad N_2(j) = \sum_{v=1}^n I(Y_v \leq Y_j, Z_v > Z_j),$$

$$N_3(j) = \sum_{v=1}^n I(Y_v > Y_j, Z_v \leq Z_j), \quad N_4(j) = \sum_{v=1}^n I(Y_v > Y_j, Z_v > Z_j),$$

it turns out that

$$T_n = \frac{1}{n^4} \sum_{j=1}^n (N_1(j)N_4(j) - N_2(j)N_3(j))^2.$$

In what follows, let $n \geq 2$. We introduce the bivariate random vectors

$$(Y_{ij}, Z_{ij}) = (f(X_i, X_j), g(X_i, X_j)), \quad 1 \leq i, j \leq n, i \neq j, \tag{1.2}$$

and adopt the above approach replacing Y_j and Z_j by Y_{ij} and Z_{ij} , respectively. The resulting test statistic is

$$\mathcal{H}\mathcal{B}\mathcal{K}\mathcal{R}_n = n \int (H_n(y, z) - F_n(y)G_n(z))^2 dH_n(y, z),$$

where now

$$H_n(y, z) = \frac{1}{n(n-1)} \sum_{\substack{i,j=1 \\ i \neq j}}^n I(Y_{ij} \leq y, Z_{ij} \leq z), \quad (y, z) \in \mathbb{R}^2,$$

and

$$F_n(y) = H_n(y, \infty), \quad G_n(z) = H_n(\infty, z), \quad y, z \in \mathbb{R},$$

are empirical distribution functions of U -statistics structure. There is an alternative expression as before,

$$\mathcal{H}\mathcal{B}\mathcal{K}\mathcal{R}_n = \frac{n}{(n(n-1))^5} \sum_{\substack{\mu, \nu=1 \\ \mu \neq \nu}}^n (N_1(\mu, \nu)N_4(\mu, \nu) - N_2(\mu, \nu)N_3(\mu, \nu))^2$$

with obvious meaning of $N_i(\mu, \nu)$, $i = 1, 2, 3, 4$, i.e.,

$$N_1(\mu, \nu) = \sum_{\substack{i,j=1 \\ i \neq j}}^n I(Y_{ij} \leq Y_{\mu\nu}, Z_{ij} \leq Z_{\mu\nu}),$$

$$N_2(\mu, \nu) = \sum_{\substack{i,j=1 \\ i \neq j}}^n I(Y_{ij} \leq Y_{\mu\nu}, Z_{ij} > Z_{\mu\nu}),$$

Download English Version:

<https://daneshyari.com/en/article/1145425>

Download Persian Version:

<https://daneshyari.com/article/1145425>

[Daneshyari.com](https://daneshyari.com)