*Mathematics*

# On Testing Hypotheses of Equality Distribution Densities

## Petre Babilua*, Elizbar Nadaraya**, Grigol Sokhadze*

*  *Department of Mathematics, Faculty of Exact and Natural Sciences, Ivane Javakhishvili Tbilisi State University, Tbilisi, Georgia*
** *Academy Member, Department of Mathematics, Faculty of Exact and Natural Sciences, Ivane Javakhishvili Tbilisi State University, Tbilisi, Georgia*

**ABSTRACT. In the paper the test of homogeneity and goodness-of-fit for checking the hypotheses of equality distribution densities is constructed. The power asymptotics of the constructed test of homogeneity and goodness-of-fit for certain types of close alternatives is also studied.** © *2016 Bull. Georg. Natl. Acad. Sci.*

**Key words:** test of homogeneity, goodness-of-fit test, power of the test, Wiener process, close alternatives

Let $X^{(i)} = \left( X_1^{(i)}, \ldots, X_{n_i}^{(i)} \right)$, $i = 1, \ldots, p$, be independent samples with sizes $n_1, n_2, \ldots, n_p$, from $p \geq 2$ general population with probability densities $f_1(x), \ldots, f_p(x)$ and it is required to test two hypotheses based on samples $X^{(i)}$, $i = 1, \ldots, p$: test of homogeneity

$$H_0: \quad f_1(x) = \cdots = f_p(x) \tag{1}$$

and goodness-of-fit test

$$H_0': \quad f_1(x) = \cdots = f_p(x) = f_0(x), \tag{2}$$

where $f_0(x)$ is a fully defined density function. In case of hypothesis $H_0$ common density function $f_0(x)$ is unknown.

In this paper the test is constructed for checking hypothesis $H_0$ and $H_0'$ against a sequence of "close" alternatives [1, 2]:

$$H_1: \quad f_i(x) = f_0(x) + \mathsf{r}(n_0) \mathfrak{c}_i \left( \frac{x - \ell_i}{\mathsf{x}(n_0)} \right) \quad \left( \mathsf{r}(n_0), \mathsf{x}(n_0) \to 0 \right),$$

$$\int \mathfrak{c}_i(x)\, dx = 0, \quad n_0 = \min(n_1, \ldots, n_p) \to \infty.$$

We consider criteria for testing $H_0$ and $H_0'$, based on statistics

$$T\left(n_1, n_2, \ldots, n_p\right) = \sum_{i=1}^{p} N_i \int \left[ \hat{f}_i(x) - \frac{1}{N} \sum_{j=1}^{p} N_j \hat{f}_j(x) \right]^2 r(x)\, dx, \tag{3}$$

where $\hat{f}_i(x)$ is Rosenblatt-Parsen kernel estimator of the distribution density $f_i(x)$:

$$\hat{f}_i(x) = \frac{a_i}{n_i} \sum_{j=1}^{n_i} K\left(a_i\left(x - X_j^{(i)}\right)\right), \qquad N_i = \frac{a_i}{n_i}, \quad N = N_i + \cdots + N_p.$$

Particular case $p = 2$ was discussed in papers [3, 4]. In this case statistics $T$ becomes more visual:

$$T\left(n_1, n_2\right) = \frac{N_1 N_2}{N_1 + N_2} \int \left(\hat{f}_1(x) - \hat{f}_2(x)\right)^2 r(x)\, dx.$$

In this paper the found limiting distribution of statistics (3) was found for hypothesis $H_1$ in case, where $n_i$ unlimited is increasing so that $n_i = n k_i$, where $n \to \infty$, and $k_i$ are constant. Let $a_1 = a_2 = \cdots = a_p = a_n$, so $a_n \to \infty$ for $n \to \infty$.

For getting limiting distribution of functional $T_n = T\left(n_1, \ldots, n_p\right)$ let us introduce conditions for functions $K(x)$, $f_0(x)$, $\xi_i(x)$, $i = 1, \ldots, p$ and $r(x)$:

(i) $K(x) \geq 0$ – function with bounded variation,

$$\int K(x)\, dx = 1, \quad x^2 K(x) \in L_1(-\infty, \infty);$$

(ii) density function $f_0(x)$ is bounded and positive on $(-\infty, \infty)$ or is bounded and positive on some finite interval $[c, d]$. Besides, it has bounded derivative in the field where it is positive;

(iii) functions $\xi_j(x)$, $j = 1, \ldots, p$, are bounded and have bounded first order derivatives, also $\xi_i(x)$ and $\xi_i^{(1)} \in L_1(-\infty, \infty)$.

(iv) weighed function $r(x)$ is piece-continuous, bounded and integrable, besides $r(\ell_k) \neq 0$, $k = 1, \ldots, p$, where $\ell_k$ is some fixed points of continuity of $r(x)$.

The following is true:

**Theorem 1.** *Le us fulfill the conditions* (i)–(iv), *also* $f_i(x) \geq 0$, $x \in (-\infty, \infty)$. *If*

$$n a_n^{-1/2} r_n^2 \chi_n \to c_0 \neq 0, \quad a_n \chi_n \to \infty, \quad r_n \chi_n = o\left(n^{-1/2}\right) \quad \left(r_n = r(n_0), \ \chi_n = \chi(n_0)\right),$$

$n a_n^{-2} \to \infty$ and $a_n^2 r_n \chi_n \to \infty$, then random variable $a_n^{1/2}\left(T_n - \nu\right)$ under hypothesis $H_1$ has normal limit distribution $\left(A(\xi), \tau^2\right)$, where

$$A(\xi) = c_0 \sum_{i=1}^{p} \left(k_i - \frac{k_i^2}{\bar{k}}\right) r(\ell_i) \int \xi_i^2(x)\, dx,$$

$$\tau^2 = 2(p-1) \int f_0^2(x) r(x)\, dx\, R(K_0), \quad K_0 = K * K,$$

$$\nu = (p-1) \int f(x) r(x)\, dx\, R(K), \quad R(g) = \int g^2(x)\, dx,$$

$$\bar{k} = k_1 + \cdots + k_p, \quad p \geq 2.$$

Conditions of Theorem 1 about $a_n$, $r_n$ and $\chi_n$ are fulfilled, for example, if: $a_n = n^u$, $r_n = n^{-r}$, $\chi_n = n^{-s}$ at $\frac{u}{2} = 1 - 2r - s$, $r + s > \frac{1}{2}$, $0 < u < \frac{1}{2}$, $0 < s < u$, and conditions about $r$, $s$ and $u$ are fulfilled, for example, if

$$u = \frac{1}{4}, \quad s = \frac{1}{5}, \quad r = \frac{27}{80}; \quad u = \frac{2}{9}, \quad s = \frac{1}{6}, \quad r = \frac{13}{36};$$

$$u = \frac{1}{5}, \quad s = \frac{1}{6}, \quad r = \frac{1}{30} \text{ etc.}$$

From Theorem 1 we will state two Corollaries:

**Corollary1.** *Let the conditions* (i), (ii) *and* (iv) *be fulfilled under* $K(x)$, $f_0(x)$ $r(x)$. *If* $na_n^{-2} \to \infty$, *then random variable* $a_n^{1/2}(T_n - \sim)$ *under hypothesis* $H_0'$ *has a normal limit distribution* $\left(0, \dagger^2\right)$.

By Corollary1 a test for checking hypothesis $H_0'$ can be constructed; critical region for checking hypothesis can be defined by inequality

$$T_n \geq d_n(r), \tag{4}$$

where

$$d_n(r) = \sim + a_n^{-1/2}\dagger\}_r,$$

$\}_r$ is the quantile of level $1 - r$ $(0 < r < 1)$ of the standard normal distribution $\Phi(x)$.

**Corollary 2.** *Under conditions of Theorem* 1 *local behavior of the power* $P_{H_1}\left(T_n \geq d_n(r)\right)$ *is as follows*

$$P_{H_1}\left(T_n \geq d_n(r)\right) \to 1 - \Phi\left(\}_r - \frac{A(\{)}{\dagger}\right),$$

*when* $n \to \infty$.

Let introduce

$$f_n^*(x) = \frac{1}{k} \sum_{j=1}^{p} k_j \mathcal{f}_j(x),$$

$$\bar{\sim}_n = \int f_n^*(x) r(x) \, dx,$$

$$\Delta_n^2 = \frac{1}{k} \sum_{i=1}^{p} k_i \Delta_{in}^2, \qquad \Delta_{in}^2 = \int \mathcal{f}_i^2(x) r(x) \, dx.$$

The Theorem is true.

**Theorem 2.** *Let all the conditions of Theorem* 1 *be fulfilled. Then*

$$a_n^{1/2}(T_n - \sim_n)\dagger_n^{-1}$$

under hypothesis $H_1$ has a normal limit distribution $\left(A(\{)\dagger^{-1}, 1\right)$, where

$$\sim_n = (p-1)R(K)\bar{\sim}_n, \qquad \dagger_n^2 = 2(p-1)R(K_0)\Delta_n^2.$$

**Proof.** It is obvious

$$a_n^{1/2}(T_n - \sim_n)\dagger_n^{-1} = a_n^{1/2}(T_n - \sim)\dagger^{-1}\left(\dagger\dagger_n^{-1}\right) + a_n^{1/2}(\sim - \sim_n)\dagger_n^{-1}.$$

As it is enough to show

$$a_n^{1/2}\left(\bar{\sim}_n - \int f_0(x)r(x)\,dx\right) = o_p(1) \tag{5}$$

and

$$\Delta_n^2 - \int f_0^2(x) r^2(x)\, dx = o_p(1).\qquad(6)$$

But (6) follows from Theorem 2.1 Bhattacharyya G. K., Roussas G. G. [5] (see also [6], [2]).

Let us proove (5). We have

$$a_n^{1/2} E\left|\int f_n^*(x) r(x)\, dx - \int f_0(x) r(x)\, dx\right| \le$$

$$\le a_n^{1/2} E\left|\int\left(f_n^*(x) - Ef_n^*(x)\right) r(x)\, dx\right| + a_n^{1/2}\int\left|Ef_n^*(x) - f_0(x)\right| r(x)\, dx =$$

$$= A_{1n} + A_{2n}.$$

It is not difficult to show

$$Ef_i^€(x) = f_0(x) + O\left(\frac{1}{a_n}\right) + \mathsf{r}_n\int K(t)\{_i\left(\frac{x-\ell_i}{\mathsf{x}_n} - \frac{t}{a_n\mathsf{x}_n}\right)dt,$$

$O(\,\cdot\,)$ uniformly in $x\in(-\infty,\infty)$. So

$$Ef_n^*(x) = f_0(x) + O\left(\frac{1}{a_n}\right) + \mathsf{r}_n\frac{1}{k}\sum_{j=1}^{p} k_j\int K(t)\{_j\left(\frac{x-\ell_j}{\mathsf{x}_n} - \frac{t}{a_n\mathsf{x}_n}\right)dt$$

it follows,

$$A_{2n} \le c_1 a_n^{-1/2} + c_2 a_n^{1/2}\mathsf{r}_n\mathsf{x}_n.$$

Next, we have

$$A_{1n} \le a_n^{1/2} E^{1/2}\left(\int\left(f_n^*(x) - Ef_n^*(x)\right) r(x)\, dx\right)^2 \le$$

$$\le c_3 a_n^{1/2}\max_{1\le j\le p}\left\{\frac{1}{n}\int f_j(u)\, du\left(\int K(t) r\left(u - \frac{t}{a_n}\right)dt\right)^2\right\}^{1/2} \le c_4\left(\frac{a_n}{n}\right)^{1/2}.$$

It follows,

$$A_{1n} + A_{2n} \le c_4\left(a_n^{-1/2} + \sqrt{a_n}\,\mathsf{r}_n\mathsf{x}_n + \left(\frac{a_n}{n}\right)^{1/2}\right) \to 0$$

as $\sqrt{a_n}\,\mathsf{r}_n\mathsf{x}_n \le a_n^2\mathsf{r}_n\mathsf{x}_n \to 0$ and $\dfrac{a_n}{n}\to 0$.

From Theorem 2 we will state two corollaries.

**Corollary 3.** *Random variable*

$$a_n^{1/2}\left(T_n - \sim_n\right)\dagger_n^{-1}$$

under hypothesis $H_0$ has a normal limit distribution $(0,1)$.

This result can be used for constructing an asymptotic test for checking hypothesis $H_0:\ f_1(x) = \cdots = f_p(x)$ (test of homogeneity); critical region can be defined by inequality:

$$T_n \ge \tilde{d}_n(\mathsf{r}) = \sim_n + a_n^{-1/2}\dagger_n\}_\mathsf{r},\qquad(7)$$

where $\}_\mathsf{r}$ is the quantile of level $1-\mathsf{r}$ of the standard normal distribution $\Phi(x)$.

**Corollary 3.** *Under conditions of Theorem2 local behavior of the power* $P_{H_1}\left(T_n \ge \tilde{d}_n(\mathsf{r})\right)$ *as follows*

$$P_{H_1}\left(T_n \geq \tilde{d}_n\left(\mathsf{r}\right)\right) \to 1 - \Phi\left(\}_\mathsf{r} - A(\{)\dagger^{-1}\right),$$

when $n \to \infty$.

**Remark 1.** *Under hypothesis* $H_1$ we have

$$F_i\left(x\right) = F_0\left(x\right) + \mathsf{r}_n \mathsf{x}_n U_i\left(\frac{x - \ell_i}{\mathsf{x}_n}\right), \quad U_i\left(u\right) = \int_{-\infty}^{u} \{_i\left(x\right) dx$$

and according to Theorem 1 $\mathsf{r}_n \mathsf{x}_n = o\left(\dfrac{1}{\sqrt{n}}\right)$. So it can be written

$$\sup_x \left|F_i\left(x\right) - F_0\left(x\right)\right| = o\left(\frac{1}{\sqrt{n}}\right). \tag{8}$$

It is well-known that the test based on deviation between empirical distribution functions, for example, criterion Kolmogorov-Smirnov and test Cramer-Mises-Smirnov (analogues to these criteria for $p \geq 2$ was constructed by **Kiefer, J** [7]) differs local close alternatives from null hypothesis, if $F_i\left(x\right) - F_0\left(x\right) = O\left(\dfrac{1}{\sqrt{n}}\right)$ uniformly in $x \in \left(-\infty, \infty\right)$, in case of (8) the above mentioned test cannot asymptotically distinguish such hypotheses from null hypothesis (limiting power will be equal with level of the test). However, tests (4) and (7) based on estimators of distribution density are more powerful asymptotically (under hypothesis $H_1$) than tests based on empirical distribution functions (analogues questions for one sample considered in paper of Rosenblatt [1]).

**Remark 2.** Tests (4) and (7) for checking hypotheses $H_0'$ and $H_0$, against alternatives $H_1$ are asymptotically strictly unbiased as $A(\{) > 0$ and equal to $0$ if and only if $\{_i\left(x\right) = 0$, almost everywhere, $i = 1, \ldots, p$.

*მათემატიკა*

# განაწილების სიმკვრივეთა ტოლობის ჰიპოთეზათა შემოწმების შესახებ

## პ. ბაბილუა\*, ე. ნადარაია\*\*, გ. სოხაძე\*

\* *ივანე ჯავახიშვილის სახ. თბილისის სახელმწიფო უნივერსიტეტი, ზუსტ და საბუნებისმეტყველო მეცნიერებათა ფაკულტეტი, მათემატიკის დეპარტამენტი, თბილისი, საქართველო*
\*\* *აკადემიის წევრი, ივანე ჯავახიშვილის სახ. თბილისის სახელმწიფო უნივერსიტეტი, ზუსტ და საბუნებისმეტყველო მეცნიერებათა ფაკულტეტი, მათემატიკის დეპარტამენტი, თბილისი, საქართველო*

ნაშრომში აგებულია ერთგვაროვნების და თანხმობის ჰიპოთეზათა შემოწმების კრიტერიუმები. მოძებნილია აგებული კრიტერიუმების ზღვარითი სიმძლავრე დაახლოებადი ალტერნატივებისთვის.

## REFERENCES

1. *Rosenblatt M.* (1975) Ann. Statist. **3**: 1-14.
2. *Nadaraya E. A.* (1989) Nonparametric estimation of probability densities and regression curves. Mathematics and its Applications (Soviet Series), 20. Kluwer Academic Publishers Group, Dordrecht.
3. *Nadaraya E. A.* (1975) Soobshch. Akad. Nauk Gruz. SSR **78**: 25-28 (in Russian).
4. *Anderson N. H., Hall P.,. Titterington D. M* (1994) J. Multivariate Anal., **50**, 1: 41-54.
5. *Bhattacharyya G. K., Roussas G. G.* (1970) Skand. Aktuarietidskr. **1969**: 201-206.
6. *Mason D. M., Nadaraya E. A., Sokhadze G. A.* (2010) Integral functionals of the density. Inst. Math. Statist. Collect. 7, Inst. Math. Stat. Beachwood, OH. 153-168.
7. *Kiefer J.* (1959) Ann. Math. Statist. **30**: 420-447.