

UDC 519.21

ROMAN I. ANDRUSHKIW, DMITRY D. KLYUSHIN, AND YURIY I. PETUNIN

A NEW TEST FOR UNIMODALITY

A distribution function (d.f.) of a random variable is unimodal if there exists a number such that d.f. is convex left from this number and is concave right from this number. This number is called a mode of d.f. Since one may have more than one mode, a mode is not necessarily unique. The purpose of this paper is to construct nonparametric tests for the unimodality of d.f. based on a sample obtained from the general population of values of the random variable by simple sampling. The tests proposed are significance tests such that the unimodality of d.f. can be guaranteed with some probability (confidence level).

1. INTRODUCTION

Testing the unimodality of a distribution function is a widely investigated issue. The most popular tests include the DIP test proposed by J.A.Hartigan and the kernel density estimation test proposed by B.W.Silverman [1–4]. However, all of these tests are computationally quite complex and asymptotic. That is why it is useful to develop elementary tests which are based on simple computational procedures and are non-asymptotic.

According to A.Ya.Khinchin, a distribution function $F(u)$ of a random variable x is unimodal if there exists a number M such that d.f. $F(u)$ is convex in $(-\infty, M)$ and concave in (M, ∞) . The number M is said to be a mode of d.f. $F(u)$. A mode can not be unique, since d.f. $F(u)$ can have several modes. Also, d.f. $F(u)$ can have break at M and be continuous in $(-\infty, M)$ and (M, ∞) .

The purpose of the paper is to construct nonparametric tests for the unimodality of d.f. $F(u)$ based on a sample x_1, x_2, \dots, x_n obtained by the simple sampling from the general population of values of a random variable x . The tests proposed are significance tests, so the unimodality of d.f. $F(u)$ can be guaranteed with some probability α (confidence level), where $\beta = 1 - \alpha$ is the significance level of a test. To formulate the tests, we introduce new estimations of the probability density (d.p.) and the distribution function based on a sample x_1, x_2, \dots, x_n .

2. UNIFIED HISTOGRAM AND MODIFIED EMPIRICAL D.F.

Let x_1, x_2, \dots, x_n be a sample obtained from a general population $F(u)$ by the simple sampling which has d.p. $F(u)$. Since these functions are unknown, we call them hypothetical. To estimate d.p., we use the relation

$$(1) \quad p(x_{n+1} \in (x_{(i)}, x_{(i+1)})) = \frac{1}{n+1},$$

where $x_{(i)}$ is the order statistics ($i = 1, 2, \dots, n$). Using estimation (1), we can define the estimation $h_n(u)$ for a hypothetical d.p.

$$h_n(u) = \begin{cases} \frac{1}{(n-1)(x_{(i+1)} - x_{(i)})}, & \text{if } u \in (x_i, x_{i+1}), \\ 0, & \text{otherwise.} \end{cases}$$

2000 *AMS Mathematics Subject Classification.* Primary 62G05.

Key words and phrases. Unimodality, distribution function, significance test.

In such a case, the probability that the value of a random variable \tilde{x} with d.p. (2) belongs to $[x_{(i)}, x_{(i+1)})$ is equal to

$$p(\tilde{x} \in (x_{(i)}, x_{(i+1)})) = \frac{1}{n+1},$$

where x_i are considered as constants. For large n , this probability is close to probability (1), so we refer to the value $h_n(u)$ as a unified histogram constructed on x_1, x_2, \dots, x_n . This histogram has some advantage over all other histograms, because it is unambiguously defined by the sample x_1, x_2, \dots, x_n . Also, the integral

$$(2) \quad \tilde{F}_n^*(u) = \int_{x_{(1)}}^u h_n(v) dv = \frac{u + (i-1)x_{(i+1)} - ix_{(i)}}{(n-1)(x_{(i+1)} - x_{(i)})}$$

is a linear spline $x_{(i)} \leq v < x_{(i+1)}$ which is a more precise estimation of the hypothetical d.f $F(u)$ than that of a piecewise empirical d.f.

$$F_n^*(u) = \frac{i}{n}, \text{ if } x_{(i)} \leq v < x_{(i+1)}.$$

We refer to the function $F_n^*(u)$ as e.d.f. and to $\tilde{F}_n^*(u)$ defined by (2) as a modified e.d.f. (m.e.d.f.). Its advantages over the conventional e.d.f. are obvious: 1) when d.f. $F(u)$ is continuous, linear splines are more precise approximations than the piecewise e.d.f $F_n^*(u)$, and 2) $\tilde{F}_n^*(u)$ is continuous, so it is possible to estimate quantiles of any order and to construct an inverse d.f. (Quetelet curve), whereas it is impossible to do by using the piecewise e.d.f. $F_n^*(u)$. However, at large n , e.d.f. $F_n^*(u)$ and $\tilde{F}_n^*(u)$ are close. Let us prove that

$$(3) \quad \left| F_n^*(u) - \tilde{F}_n^*(u) \right| \leq \frac{1}{n}.$$

Indeed, for all $u \in [x_i, x_{i+1})$, the following relation holds:

$$\left| F_n^*(u) - \tilde{F}_n^*(u) \right| = \left| \frac{u + (i-1)x_{(i+1)} - ix_{(i)}}{n(n-1)(x_{(i+1)} - x_{(i)})} - \frac{i}{n} \right| = \left| \frac{nu - (n-i)x_{(i+1)} + ix_{(i)}}{n(n-1)(x_{(i+1)} - x_{(i)})} \right|.$$

Granting that

$$\tilde{F}_n^*(u) = \frac{u - x_{(i+1)} + i(x_{(i+1)} - x_{(i)})}{(n-1)(x_{(i+1)} - x_{(i)})} = \frac{1}{n-1} \left[\frac{u - x_{(i+1)}}{x_{(i+1)} - x_{(i+1)}} + 1 \right],$$

we have

$$\frac{i-1}{n-1} \leq \tilde{F}_n^*(u) \leq \frac{i}{n-1}$$

and

$$\begin{aligned} \frac{i-1}{n-1} - F_n^*(u) &\leq \tilde{F}_n^*(u) - F_n^*(u) \leq \frac{i}{n-1} - F_n^*(u), \\ \frac{n(i-1) - (n-1)i}{(n-1)i} &\leq \tilde{F}_n^*(u) - F_n^*(u) \leq \frac{ni - (n-1)i}{n(n-1)}, \\ \frac{i-n}{n(n-1)} &\leq \tilde{F}_n^*(u) - F_n^*(u) \leq \frac{1}{n}; \end{aligned}$$

hence,

$$-\frac{1}{n} \leq \tilde{F}_n^*(u) - F_n^*(u) \leq \frac{1}{n},$$

i.e.

$$(4) \quad \left| \tilde{F}_n^*(u) - F_n^*(u) \right| \leq \frac{1}{n}.$$

Estimation (4) implies that m.e.d.f has similar asymptotic properties as conventional e.d.f., i.e. it is consistent, asymptotically unbiased, etc.

3. CONFIDENCE LIMITS FOR HYPOTHETICAL D.F.

Let us define the lower and upper bounds of a hypothetical d.f. $F(u)$ by means of an empirical d.f., under the assumption that $F(u)$ is continuous and strictly increasing. This problem was solved in [5–10]. Hence, given a significance level β^* (e.g., $\beta^* = 0.05$), we can define ε so that

$$p\left(\Delta = \max_{x_{(1)} \leq u \leq x_{(n)}} |F(u) - F_n(u)| > \varepsilon\right) = \beta^*.$$

It follows that, for a given β^* , we can find ε according to statistical tables [11] and construct a strip Π_{β^*} , whose bounds are stepwise linear: $y = F_n^*(u) + \varepsilon$ and $y = F_n^*(u) - \varepsilon$. The strip Π_{β^*} completely covers the true d.f. $y = F(u)$ with the confidence probability $\alpha^* = 1 - \beta^*$. Hereinafter, we refer to the strip Π_{β^*} as the confidence strip for d.f. with significance level β^* constructed for the empirical d.f.

4. TEST FOR UNIMODALITY BASED ON E.D.F.

Let x_1, x_2, \dots, x_n be a sample obtained from a general population G by the simple sampling with continuous and strictly monotone d.f. $F(u)$. Using this sample, we construct the empirical d.f. $F_n^*(u)$ and the strip Π_{β^*} . Denote, by $\varphi(u)$, the upper bound of Π_{β^*} described by the equation $y = F_n^*(u) + \varepsilon$, and let $\psi(u)$ be the lower bound of Π_{β^*} described by the equation $y = F_n^*(u) - \varepsilon$. Then

$$p(\varphi(u) \leq F(u) \leq \psi(u)) = \alpha^* = 1 - \beta^*.$$

Definition 1. Let $y = \varphi(u)$ be an arbitrary function defined on $[a, b]$. Then the set

$$G_U = \{(u, y) : y \geq \varphi(u), a \leq u \leq b\}$$

is an epigraph of $\varphi(u)$, and the set $G_L = \{(u, y) : y \leq \varphi(u), a \leq u \leq b\}$ is a subgraph of $\varphi(u)$.

Definition 2. The lower bound of a convex hull of the epigraph of a function $\varphi(u)$ is a convex minorant of $\varphi(u)$,

$$\varphi_{\text{inf}}(u) = \inf \left\{ v : (u, v) \in \text{conv}_{a \leq u \leq b} G_U \right\},$$

where $\text{conv } G_U$ is the convex hull of G_U . Analogously, the upper bound of a convex hull of the subgraph of a function $\varphi(u)$ is a concave majorant of $\varphi(u)$:

$$\psi_{\text{sup}}(u) = \sup \left\{ v : (u, v) \in \text{conv}_{a \leq u \leq b} G_L \right\}.$$

Theorem 1. Let $\varphi_{\text{inf}}(u)$ and $\psi_{\text{sup}}(u)$ be the convex minorant and concave majorant of $\varphi(u)$ and $\psi(u)$, respectively, and

$$c = \sup \{ u : \varphi_{\text{inf}}(u) \leq \psi(u), x_{(1)} \leq u \leq x_{(n)} \},$$

$$d = \inf \{ u : \psi_{\text{sup}}(u) \geq \varphi(u), x_{(1)} \leq u \leq x_{(n)} \}.$$

Then, the hypothetical distribution $F(u)$ is unimodal iff

1) $\varphi_{\text{inf}}(u) \geq \psi(u)$ or $\psi_{\text{sup}}(u) \leq \varphi(u) \quad \forall u \in [x_{(1)}, x_{(n)}]$;

or

2) $c \geq d$.

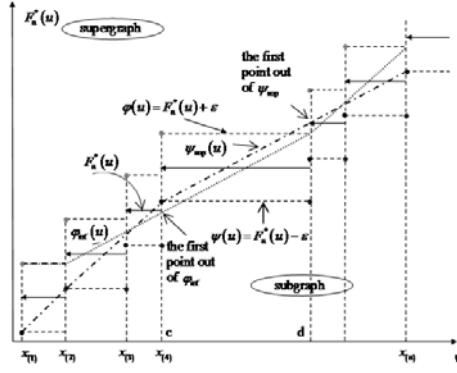


Fig. 1. If $c < d$, the unimodality is absent.

Moreover, the significance level of this criterion is β^* .

Proof. Necessity. Suppose that the hypothetical d.f. $F(u)$ is unimodal, and M is its mode. If $M \leq x_{(1)}$ or $M \geq x_{(n)}$, then $F(u)$ on $[x_{(1)}, x_{(n)}]$ can be convex or concave. Then condition 1) holds.

If $x_{(1)} \geq M \leq x_{(n)}$, then $F(u)$ is convex on $[x_{(1)}, M]$ and concave on $[M, x_{(n)}]$. In such a case, it follows from Definition 2 (see Fig. 1) that $\varphi_{\text{inf}}(u) \geq F(u)$ on $[x_{(1)}, M]$. Also, $F(u) \geq \psi_{\text{inf}}(u) \forall u \in [x_{(1)}, M]$, so $\varphi_{\text{inf}}(u) \geq \psi(u)$, and $d \geq M$.

On the other hand, $F(u) \geq \psi_{\text{sup}}(u)$, as far as $F(u)$ is concave on $[M, x_{(n)}]$. Definition 2 implies that $F(u) \geq \psi_{\text{sup}}(u)$ on $[M, x_{(n)}]$. Also, $\varphi(u) \geq F(u) \forall u \in [M, x_{(n)}]$. Thus, $\varphi(u) \geq \psi_{\text{sup}}(u) \forall u \in [M, x_{(n)}]$ and $d \leq M$. Consequently, $c \geq d$, and condition 2) holds.

Sufficiency. Note that $\varphi(u)$ and $\psi(u)$ are increasing. If condition 1) holds, then

$$\psi(u) \leq \varphi_{\text{sup}}(u) \leq \varphi(u) \quad \forall u \in [x_{(1)}, x_{(n)}]$$

or

$$\psi(u) \leq \psi_{\text{inf}}(u) \leq \varphi(u) \quad \forall u \in [x_{(1)}, x_{(n)}].$$

Thus, $\varphi_{\text{sup}}(u)$ (or $\psi_{\text{inf}}(u)$) lies in the strip Π_β . Therefore, $\varphi_{\text{sup}}(u)$ (or $\psi_{\text{inf}}(u)$) can be used as an estimation of the hypothetical d.f. $F(u)$ of a general population G . Since $F(u) = \varphi_{\text{inf}}(u)$ or $F(u) = \psi_{\text{inf}}(u)$, the hypothetical d.f. increases, is convex or concave on $[x_{(1)}, x_{(n)}]$, and is unimodal. The significance level of this test is β^* .

Now, we suppose that condition 2) holds, i.e. $c \geq d$. Put $\hat{F}(u) = \varphi_{\text{inf}}(u)$, if $u \in [x_{(1)}, c]$, and $\hat{F}(u) = \psi_{\text{sup}}(u)$, if $u \in (c, x_{(n)})$. It is easy to see that $\hat{F}(u)$ lies in Π_β , because $c \geq d$. Also, $\hat{F}(u)$ is convex on $[x_{(1)}, c]$ and concave on $(c, x_{(n)})$. Let us prove that $\hat{F}(u) \geq \hat{F}(c+0) = \lim_{u \rightarrow c, u > c} \gamma = \gamma$. Indeed, if $\gamma < \hat{F}(c)$, then $\gamma \notin \Pi_\beta$. Therefore, the abscissa of the first exit point d , where $\varphi_{\text{inf}}(u)$ exceeds the bounds Π_β while moving from $x_{(n)}$ to $x_{(1)}$, is greater than c . This contradicts condition 2. Thus, $\hat{F}(u)$ increases, is convex on $[x_{(1)}, c]$, and concave on $(c, x_{(n)})$. But $\hat{F}(u)$ can have a breakpoint in c . To exclude this breakpoint, we take $\varepsilon > 0$ sufficiently small so that the segment with the ends $(c - \varepsilon, \hat{F}(c - \varepsilon))$ completely lies in Π_β . Then the function

$$\hat{F}_\varepsilon(u) = \begin{cases} \hat{F}(u), & \text{if } u \in [x_{(1)}, c - \varepsilon], \\ \hat{u} \frac{\gamma - \hat{F}(c - \varepsilon)}{\varepsilon} + \gamma - c \frac{\gamma - \hat{F}(c - \varepsilon)}{\varepsilon}, & \text{if } u \in (c - \varepsilon, c), \\ \hat{F}(u), & \text{if } u \in [c, x_{(n)}], \end{cases}$$

increases, is continuous, convex on $[x_{(1)}, c]$, concave on $[c, x_{(n)}]$, and its graph lies in Π_β . Thus, d.f. $\hat{F}_\varepsilon(u)$ is unimodal, and we can consider it as an estimation of the hypothetical d.f. of a general population G . The significance level of this test is β . Theorem 1 is proved.

Remark 1. Theorem 1 has the following geometric sense: let c be the abscissa of the first exit point, where the convex minorant $\varphi_{inf}(u)$ exceeds the upper bound of Π_β while moving from the maximal order statistics to the minimal one, and let d be the abscissa of the first exit point, where the convex minorant $\psi_{sup}(u)$ exceeds the upper bound of Π_β while moving from the minimal order statistics to the maximal one. Then the hypothetical d.f. $F(u)$ is unimodal iff the exit points c and d lie outside $[x_{(1)}, x_{(n)}]$ or $c \geq d$.

5. TEST FOR UNIMODALITY BASED ON M.E.D.F.

The confidence strip $\hat{\Pi}_\beta$ for a hypothetical d.f. can be constructed on m.e.d.f. $\tilde{F}_n^*(u)$ in the following way: let the significance level β^* be given, let ε be the width of Π_β , and let

$$p\left(\Delta = \max_{x_{(1)} \leq u \leq x_{(n)}} |F(u) - F_n^*(u)| > \varepsilon\right) = \beta^*.$$

We put $\tilde{\varphi}(u) = \tilde{F}_n^*(u) + \varepsilon + \frac{1}{n}$ and $\tilde{\psi}(u) = \tilde{F}_n^*(u) - \varepsilon - \frac{1}{n}$. It is easy to see that $\tilde{\Pi}_{\beta^*}$ with lower bound $\tilde{\psi}(u)$ and upper bound $\tilde{\varphi}(u)$ has the significance level not exceeding β^* . Indeed, by the virtue of (4),

$$\left|F(u) - \tilde{F}_n^*(u)\right| = \left|F(u) - F_n^*(u) + F_n^*(u) - \tilde{F}_n^*(u)\right| \leq |F(u) - F_n^*(u)| + \frac{1}{n}$$

Therefore,

$$\tilde{\Delta} = \max_{x_{(1)} \leq u \leq x_{(n)}} |F(u) - F_n^*(u)| < \Delta + \frac{1}{n}$$

Hence,

$$\begin{aligned} p\left(\tilde{\Delta} - \frac{1}{n} \geq \varepsilon\right) &\leq p(\Delta \geq \varepsilon) = \beta^*, \\ p\left(\tilde{\Delta} \geq \varepsilon + \frac{1}{n}\right) &\leq p(\tilde{\Delta} \geq \tilde{\varepsilon}) = \beta^*. \end{aligned}$$

Thus, the significance level of $\tilde{\Pi}_{\beta^*}$ does not exceed β^* . Since we increase the validity of the test by selecting β^* as a significance level, we can use the m.e.d.f. $\tilde{F}_n^*(u)$ to construct $\hat{\Pi}_\beta$ without decrease in the significance level. However, doing this, we increase the width of $\tilde{\Pi}_{\beta^*}$ by $\frac{1}{n}$ relative to Π_{β^*} . For moderate samples ($30 \leq n \leq 200$), this increment varies from 7 to 13

Now, we can formulate the test for the unimodality of a hypothetical d.f. based on m.e.d.f.

Theorem 2. Let $\tilde{\varphi}_{inf}(u)$ and $\tilde{\psi}_{sup}(u)$ be the convex minorant and concave majorant of $\tilde{\varphi}(u)$ and $\tilde{\psi}(u)$, respectively, and let

$$\begin{aligned} c &= \sup \left\{ u : \tilde{\varphi}_{inf}(u) \leq \tilde{\psi}(u), x_{(1)} \leq u \leq x_{(n)} \right\}, \\ d &= \inf \left\{ u : \tilde{\psi}_{sup}(u) \geq \tilde{\varphi}(u), x_{(1)} \leq u \leq x_{(n)} \right\}. \end{aligned}$$

Then, the hypothetical distribution $F(u)$ is unimodal iff

1) $\tilde{\varphi}_{inf}(u) \geq \tilde{\psi}(u)$ or $\tilde{\psi}_{sup}(u) \leq \tilde{\varphi}(u) \forall u \in [x_{(1)}, x_{(n)}]$;

or

2) $c \geq d$.

Moreover, the significance level of this criterion is β^* .

The proof of Theorem 2 is similar to that of Theorem 1.

6. CONCLUSION

It is shown in [12] that if the distribution function of a general population is unimodal, then the confidence interval $(m(x) - 3\sigma(x), m(x) + 3\sigma(x))$, where $m(x)$ is the mathematical expectation of G and $\sigma(x)$ is the standard deviation of G , has the significance level which does not exceed 0.05. That is why, this nonparametric test for unimodality can be used to construct the confidence interval for the bulk of the general population G .

BIBLIOGRAPHY

1. B.W. Silverman, *Using kernel density estimates to investigate multimodality*, J. of the Royal Statistical Society B **43** (1981), 97-99.
2. J.A. Hartigan, *Computation of the dip statistics to test for unimodality*, Applied Statistics **34** (1985), 320-325.
3. J.A. Hartigan, *The span test of multimodality*, Classification and Related Methods of Data Analysis, (H. H. Bock, ed.), North-Holland, Amsterdam, 1988, pp. 229-236.
4. J.A. Hartigan, S. Mohanty, *The RUNT test for multimodality*, Applied Statistics **9** (1992), 63-70.
5. A.N. Kolmogoroff, *Determinatione empirica di una legge di distribuzione*, Giornale Instit. Ital. Attuari **4** (1933), 83-91.
6. N.V. Smirnov, *Sur les ecart de la courbe de distribution empiric*, Mat. Sb. **6** (1939), 3-26.
7. A. Wald, J. Wolfowitz, *Confidence limits for continuous distribution functions*, Ann. Math. Statist. **10** (1939), 199-326.
8. W. Feller, *On the Kolmogorov-Smirnov limit theorems for empirical distributions*, Ann. Math. Statist. **19** (1948), 177-189.
9. F.J. Massey, *A note on the estimation of a distribution function by confidence limits*, Ann. Math. Statist. **21** (1950), 125-128.
10. Z.W. Birnbaum, F.H. Tingey, *One-sided confidence contours for distribution functions*, Ann. Math. Statist. **22** (1951), 592-596.
11. B.L. Van der Waerden, *Mathematische Statistik*, Springer, Berlin, 1957.
12. D.F. Vysochanskij, Yu.I. Petunin, *Justification of the 3- σ rule for unimodal distribution*, Theor. Probability Math. Stat. **21** (1980), 25-36.

DEPARTMENT OF MATHEMATICAL SCIENCES AND CENTER FOR APPLIED MATHEMATICS AND STATISTICS, NEW JERSEY INSTITUTE OF TECHNOLOGY, NEWARK, NJ 07102, USA

TARAS SHEVCHENKO KYIV NATIONAL UNIVERSITY, DEPARTMENT OF CYBERNETICS, 64, VOLODYMYRSKA STR., KYIV 01033, UKRAINE
E-mail: vm214@dcp.kiev.ua