

CLASSIFICATIONS WITH RELATIONS II: ASYMPTOTIC TESTING AND ESTIMATION

JAN ŘEHÁK, BLANKA ŘEHÁKOVÁ

The asymptotic statistical theory for measures derived from the D -model provides the means for analysis of generalized categorical data with relations among categories that are specified by a matrix $\mathbf{D} = \|d_{ij}\|$. This concept was introduced in [9] and developed in [10]–[19]. Independently it appeared later in the context of C.R. Rao's concept of diversity (see e.g. [4]–[7]). Weights (scores) d_{ij} express *dissimilarities* between categories i and j of the particular variable. They may be interpreted as values of a loss function as well. In this paper we present some asymptotic results for the coefficient of explanatory power (Section 2) and the coefficient of partial association (Section 3). The necessary definition and results are reviewed in Section 1. Further we deal with testing hypotheses of goodness-of-fit, homogeneity of independent samples and marginal homogeneity in square contingency tables for large samples (Section 4). Residual analysis in the homogeneity problem is also included (Section 4). Most of this paper is based on the results of [18].

1. DEFINITIONS AND PRELIMINARIES

A *generalized categorical variable* \mathbf{A} is given by a list of its values (categories) and by a matrix of scores *generating its type*:

$$(1) \quad \mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}.$$

The matrix \mathbf{D} is supposed to be a $K \times K$ real symmetric matrix with nonnegative elements d_{ij} ($d_{ii} = 0$ for all i). The more unlike the categories \mathbf{a}_i and \mathbf{a}_j are, the greater is the score d_{ij} .

A *generalized variance* of a distribution \mathbf{f} on $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}$ is defined as

$$(2) \quad \text{Gvar } \mathbf{f} = \mathbf{f}' \mathbf{D} \mathbf{f},$$

where $\mathbf{f} \in \mathbf{Q}_K$ and \mathbf{Q}_K is the set of all possible K -dimensional probability vectors.

A *distance* between distributions \mathbf{f} and \mathbf{g} of $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}$ is defined as

$$(3) \quad D(\mathbf{f}, \mathbf{g}) = [(\mathbf{f} - \mathbf{g})' \mathbf{D} (\mathbf{g} - \mathbf{f})]^{1/2}$$

whenever D is such a matrix that $(\mathbf{f} - \mathbf{g})' \mathbf{D}(\mathbf{g} - \mathbf{f}) \geq 0$ for all $\mathbf{f}, \mathbf{g} \in \mathbf{Q}_K$. This requirement is fulfilled iff $\mathbf{D}^* = \|\|d_{ij}^*\|\|$

$$(4) \quad d_{ij}^* = d_{iK} + d_{Kj} - d_{ij}$$

is a $(K - 1) \times (K - 1)$ Gramian matrix. Moreover the function $D(\cdot, \cdot)$ is a metrics (semimetrics) on \mathbf{Q}_K iff the matrix \mathbf{D}^* is a positive definite (positive semidefinite) matrix. Further it holds that

$$(5) \quad D(\mathbf{f}, \mathbf{g}) = [(\mathbf{f}^* - \mathbf{g}^*)' \mathbf{D}^*(\mathbf{f}^* - \mathbf{g}^*)]^{1/2},$$

where $\mathbf{f}^* = (f_1, \dots, f_{K-1})'$, $\mathbf{g}^* = (g_1, \dots, g_{K-1})'$.

Let $\mathbf{f}_{(1)}, \dots, \mathbf{f}_{(R)}$ be R distributions on $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}$ for the strata that are determined by values $\mathbf{b}_1, \dots, \mathbf{b}_R$ of a nominal variable \mathbf{B} , and let $\mathbf{w} \in \mathbf{Q}_R$ with positive components w_r . Considering a vector $\mathbf{f} = \sum w_r \mathbf{f}_{(r)}$ then the decompositional property of the generalized variance holds:

$$(6) \quad \begin{aligned} \text{Gvar } \mathbf{f} &= \sum_{r=1}^R w_r \text{Gvar } \mathbf{f}_{(r)} + \frac{1}{2} \sum_{r=1}^R \sum_{s=1}^R w_r w_s D^2(\mathbf{f}_{(r)}, \mathbf{f}_{(s)}) = \\ &= \sum_{r=1}^R w_r \text{Gvar } \mathbf{f}_{(r)} + \sum_{r=1}^R w_r D^2(\mathbf{f}_{(r)}, \mathbf{f}). \end{aligned}$$

A coefficient of explanatory power of decomposition $\delta_{\mathbf{A}|\mathbf{B}}$ is defined as a relative portion of the explained variability of the dependent variable \mathbf{A} by the nominal variable \mathbf{B} that generates the decomposition, i.e.

$$(7) \quad \delta_{\mathbf{A}|\mathbf{B}} = 1 - \frac{\sum_{r=1}^R w_r \text{Gvar } \mathbf{f}_{(r)}}{\text{Gvar } \mathbf{f}}$$

2. ASYMPTOTIC DISTRIBUTION OF $\delta_{\mathbf{A}|\mathbf{B}}$

The asymptotic distribution of $\hat{\delta}_{\mathbf{A}|\mathbf{B}}$ (the sample analog of $\delta_{\mathbf{A}|\mathbf{B}}$ that is its consistent estimate) has been investigated by Řeháková [18].

Consider a two-dimensional categorization (\mathbf{B}, \mathbf{A}) , where $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_R\}$ is a nominal variable and $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}$. Let $p_{rk} = P(\mathbf{B} = \mathbf{b}_r, \mathbf{A} = \mathbf{a}_k)$, $p_{rk} > 0$ for all r, k and $\sum_r \sum_k p_{rk} = 1$. Suppose that a sample of size n is drawn from a distribution with probabilities p_{rk} . Such a sample may be described by the multinomial distribution $\mathcal{M}(n; p_{11}, \dots, p_{RK})$ with RK cells. Let f_{rk} be observed relative frequencies.

Theorem 1. The asymptotic distribution of $\sqrt{(n)}(\hat{\delta}_{\mathbf{A}|\mathbf{B}} - \delta_{\mathbf{A}|\mathbf{B}})$, where

$$(8) \quad \hat{\delta}_{\mathbf{A}|\mathbf{B}} = 1 - \frac{\sum_{r=1}^R p_r \text{Gvar } \mathbf{p}_{(r)}}{\text{Gvar } \mathbf{p}},$$

$\hat{\delta}_{\mathbf{A}|\mathbf{B}} = \hat{\delta}_{\mathbf{A}|\mathbf{B}}(f_{11}, \dots, f_{RK})$ is under a random sample from the multinomial distribu-

tion $\mathcal{M}(n; p_{11}, \dots, p_{RK})$ a normal distribution $\mathcal{N}(0, \sigma^2)$ with

$$(9) \quad \sigma^2 = \frac{1}{\zeta^4} \sum_{b=1}^R \sum_{a=1}^K p_{ba} [v(2d_a^* - \zeta) - \zeta(2d_{a|b}^* - \text{Gvar } \mathbf{p}_{(b)})]^2$$

under the condition that $\sigma^2 \neq 0$. Here we denote

$$(10) \quad v = \sum_{r=1}^R p_{r+} \text{Gvar } \mathbf{p}_{(r)}, \quad \zeta = \text{Gvar } \mathbf{p},$$

$$\mathbf{p}_{(r)} = (p_{r1|p_{r+}}, \dots, p_{rK|p_{r+}})' = (p_{1|r}, \dots, p_{K|r})', \quad \mathbf{p} = (p_{+1}, \dots, p_{+K})'$$

$$d_{a|b}^* = \sum_{k=1}^K p_{k|b} d_{ka}, \quad d_a^* = \sum_{k=1}^K p_{+k} d_{ka}.$$

Further it holds that

$$\frac{\sqrt{(n)} (\delta_{A|B} - \delta_{A|B})}{\sigma(f_{11}, \dots, f_{RK})}$$

has the normal distribution $\mathcal{N}(0, 1)$. The asymptotic variance σ^2 is equal to zero if and only if all $\mathbf{p}_{(b)}$ ($b = 1, \dots, R$) belong to the same class of equivalence generated by the $D(\cdot, \cdot)$, i.e. if and only if the coefficient $\delta_{A|B}$ is equal to zero.

Proof. The form of the asymptotic distribution follows from the δ method (see [3], p. 430). For the possibility of simplifying the computation of the asymptotic variance see [2]. The details are given in [18]. To prove the last part let us note that

$$\sigma^2 = 0 \Leftrightarrow v(2d_a^* - \zeta) - \zeta(2d_{a|b}^* - \text{Gvar } \mathbf{p}_{(b)}) = F_{ba} = 0 \text{ for all } a \text{ and } b.$$

1. Let $D(\mathbf{p}_{(b)}, \mathbf{p}_{(b')}) = 0$ for all choices of b and b' . Expressing $\delta_{A|B}$ as

$$[\frac{1}{2} \sum_r \sum_s p_{r+} p_{s+} D^2(\mathbf{p}_{(r)}, \mathbf{p}_{(s)})] / \text{Gvar } \mathbf{p}$$

(see (6)) we have $\delta_{A|B} = 0$ and therefore $v = \zeta$. From the relation

$$\zeta = v + \sum_r p_{r+} D^2(\mathbf{p}_{(r)}, \mathbf{p})$$

(see (6)) it follows that under the condition $v = \zeta$ it holds that $\sum_r p_{r+} D^2(\mathbf{p}_{(r)}, \mathbf{p}) = 0$ and therefore $D^2(\mathbf{p}_{(r)}, \mathbf{p}) = 0$ for all r . Finally

$$F_{ba} = \zeta[2(d_a^* - d_{a|b}^*) - (\text{Gvar } \mathbf{p} - \text{Gvar } \mathbf{p}_{(b)})].$$

According to Corollary 1 and 2 of Theorem 4 in [15] it follows that $F_{ba} = 0$ for all a and b .

2. Conversely, let $F_{ba} = 0$ for all a and b . Then also $F_{ba_1} - F_{ba_2} = 0$. Taking the difference of this value between b_1 and b_2 we shall obtain that

$$0 = -\zeta(d_{a_1|b_1}^* - d_{a_2|b_1}^*) + \zeta(d_{a_1|b_2}^* - d_{a_2|b_2}^*)$$

and consequently

$$\sum_k p_{k|b_1} (d_{ka_1} - d_{ka_2}) = \sum_k p_{k|b_2} (d_{ka_1} - d_{ka_2}),$$

written in the vector form: $(\mathbf{d}_{a_1} - \mathbf{d}_{a_2})' \mathbf{v} = 0$, where

$$\mathbf{d}_j = (d_{1j}, \dots, d_{Kj})', \quad \mathbf{v} = \mathbf{v}(b_1, b_2) = \mathbf{p}_{(b_1)} - \mathbf{p}_{(b_2)}.$$

From here it follows that $\mathbf{d}'_1 \mathbf{v} = \mathbf{d}'_2 \mathbf{v}, \dots, \mathbf{d}'_1 \mathbf{v} = \mathbf{d}'_K \mathbf{v} \Rightarrow \mathbf{d}'_1 \mathbf{v} = 1/K(\sum_j \mathbf{d}'_j \mathbf{v}) = \mathbf{d}' \mathbf{v}$ and similarly $\mathbf{d}'_2 \mathbf{v} = \mathbf{d}' \mathbf{v}, \dots, \mathbf{d}'_K \mathbf{v} = \mathbf{d}' \mathbf{v}$, which results in $(\mathbf{d}_1 - \mathbf{d}_K)' \mathbf{v} = 0, \dots, (\mathbf{d}_{K-1} - \mathbf{d}_K)' \mathbf{v} = 0$. Written in components:

$$0 = \sum_{j=1}^K (d_{kj} - d_{Kj}) v_j \quad (k = 1, \dots, K-1).$$

If we take account of $\sum_{j=1}^K v_j = 0$ we shall obtain that

$$0 = \sum_{j=1}^{K-1} (d_{kj} - d_{Kj} - d_{kK}) v_j \quad (k = 1, \dots, K-1),$$

which rewritten in matrix notation gives $\mathbf{D}^*(\mathbf{p}_{(b_1)}^* - \mathbf{p}_{(b_2)}^*) = \mathbf{0}$, where \mathbf{D}^* is the $(K-1) \times (K-1)$ matrix with elements given by (4) and $\mathbf{p}_{(b_1)}^* = (p_{1|b_1}, \dots, p_{K-1|b_1})'$, $\mathbf{p}_{(b_2)}^*$ is defined similarly. From here we can see that $D^2(\mathbf{p}_{(b_1)}, \mathbf{p}_{(b_2)}) = 0$ for all b_1 and b_2 and therefore $\delta_{A|B} = 0$. \square

Corollary. If \mathbf{D} generates a metrics on \mathbf{Q}_K then every class of equivalence contains only one element and therefore under the condition that $p_{ba} > 0$ for all, b, a , the asymptotic variance σ^2 is equal to zero if and only if $\mathbf{p}_{(1)} = \dots = \mathbf{p}_{(R)} = \mathbf{p}$, i.e. if and only if $\delta_{A|B} = 0$.

Remark. Equivalence $\sigma^2 = 0 \Leftrightarrow \delta_{A|B} = 0$ shows that $\hat{\delta}_{A|B}$ is not convenient for testing the hypothesis $\delta_{A|B} = 0$.

Now we consider product multinomial sampling with observed conditioned relative frequencies $f_{1|r}, \dots, f_{K|r}$ ($r = 1, \dots, R$).

Theorem 2. The asymptotic distribution of $\sqrt{(n)}(\hat{\delta}_{A|B} - \delta_{A|B})$, where $\delta_{A|B}$ is given by (8), $\hat{\delta}_{A|B} = \hat{\delta}_{A|B} f_{1|1}, \dots, f_{K|R}$ is under a random sample from the R independent multinomial distributions $\mathcal{M}(n_{r+}; \mathbf{p}_{(r)})$ ($n_{r+} = n\omega_r$ is known and positive, $\sum \omega_r = 1$, $p_{a|r} > 0$ for all r and k) a normal distribution $\mathcal{N}(0, \sigma^2)$ with

$$(11) \quad \sigma^2 = \frac{4}{\zeta^4} \sum_{b=1}^R \frac{p_{b+}}{\omega_b} \sum_{a=1}^K p_{a|b} \left[v(d_a^* - \sum_{j=1}^K p_{j|b} d_j^*) - \zeta(d_{a|b}^* - \text{Gvar } \mathbf{p}_{(b)}) \right]^2,$$

where $v, \zeta, d_{a|b}^*, d_a^*$ are given by (10), under the condition that $\sigma^2 \neq 0$. Further it holds that

$$\frac{\sqrt{(n)}(\hat{\delta}_{A|B} - \delta_{A|B})}{\sigma(f_{1|1}, \dots, f_{K|R})}$$

has the normal distribution $\mathcal{N}(0, 1)$. The asymptotic variance σ^2 is equal to zero if and only if all $\mathbf{p}_{(b)}$ ($b = 1, \dots, R$) belong to the same class of equivalence generated by the $D(\cdot, \cdot)$, i.e. if and only if the coefficient $\delta_{A|B}$ is equal to zero.

The proof is analogical to that given in Theorem 1 and therefore will be omitted. It can be found in [18].

The remainder of this section is devoted to a bias of $\hat{\delta}_{\mathbf{A}|\mathbf{B}}$ under multinomial and product multinomial sampling.

Theorem 3. Under the conditions of Theorem 1 the bias of $\hat{\delta}_{\mathbf{A}|\mathbf{B}}$ is

$$(12) \quad E(\hat{\delta}_{\mathbf{A}|\mathbf{B}}) - \delta_{\mathbf{A}|\mathbf{B}} = \frac{1}{n} \left\{ \frac{v}{\zeta} + \frac{1}{\zeta} \sum_{b=1}^R \text{Gvar } \mathbf{p}_{(b)} + \frac{2}{\zeta^2} \sum_{b=1}^R \sum_{a=1}^K (2d_{a|b}^* - \text{Gvar } \mathbf{p}_{(b)}) d_a^* p_{ba} - \frac{4v}{\zeta^3} \sum_{a=1}^K (d_a^*)^2 p_{+a} \right\} + O(n^{-2}).$$

Under the condition of Theorem 2 the bias of $\hat{\delta}_{\mathbf{A}|\mathbf{B}}$ is

$$(13) \quad E(\hat{\delta}_{\mathbf{A}|\mathbf{B}}) - \delta_{\mathbf{A}|\mathbf{B}} = \frac{1}{n} \sum_{b=1}^R \frac{p_{b+}^2}{\zeta \omega_b} \left\{ \text{Gvar } \mathbf{p}_{(b)} \left(\frac{1}{p_{b+}} - \frac{v}{\zeta} \right) + \frac{4}{\zeta} \sum_{a=1}^K (d_{a|b}^* - \text{Gvar } \mathbf{p}_{(b)}) d_a^* p_{a|b} - \frac{4v}{\zeta^2} \left[\sum_{a=1}^K (d_a^*)^2 p_{a|b} - \left(\sum_{a=1}^K d_a^* p_{a|b} \right)^2 \right] \right\} + O(n^{-2}).$$

The proof that is based on the Taylor expansion is rather long and involves tedious algebraic manipulations. It is given in full details by Řeháková [18].

3. ASYMPTOTIC DISTRIBUTION OF $\hat{\delta}_{\mathbf{A}|\mathbf{C}(\mathbf{B})}$

Consider a three-dimensional categorization $(\mathbf{B}, \mathbf{C}, \mathbf{A})$, $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_R\}$, $\mathbf{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_S\}$, $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}$. The coefficient of partial association $\delta_{\mathbf{A}|\mathbf{C}(\mathbf{B})} = \delta_{(\mathbf{C}-\mathbf{A}|\mathbf{B})}$, where \mathbf{A} is a generalized categorical variable and \mathbf{B}, \mathbf{C} are supposed to be nominal ones, has been introduced by Řeháková [17] in usual way as

$$(14) \quad \delta_{\mathbf{A}|\mathbf{C}(\mathbf{B})} = \frac{\delta_{\mathbf{A}|\mathbf{B} \times \mathbf{C}} - \delta_{\mathbf{A}|\mathbf{B}}}{1 - \delta_{\mathbf{A}|\mathbf{B}}}.$$

The partial association is considered in the context of contingency tables with one dependent variable of general type and two nominal predictors. The case of a three-way contingency table does not restrict the generality since the variables \mathbf{B} and \mathbf{C} may be considered as a combination of several nominal predictors. The full range of possible applications and interpretations has been given in [14] and [16]. The asymptotic properties are summarized by the following theorem.

Theorem 4. The asymptotic distribution of $\sqrt{(n)} (\hat{\delta}_{\mathbf{A}|\mathbf{C}(\mathbf{B})} - \delta_{\mathbf{A}|\mathbf{C}(\mathbf{B})})$, where

$$(15) \quad \delta_{\mathbf{A}|\mathbf{C}(\mathbf{B})} = 1 - \frac{\sum_{r=1}^R \sum_{s=1}^S p_{rs+} \text{Gvar } \mathbf{p}_{(rs)}}{\sum_{r=1}^R p_{r+} + \text{Gvar } \mathbf{p}_{(r)}},$$

$\delta_{A|C(B)} = \delta_{A|C(B)}(f_{111}, \dots, f_{RSK})$ is under a random sample from the multinomial distribution $\mathcal{M}(n; p_{111}, \dots, p_{RSK})$ ($p_{rsk} > 0$ for all r, s, k) a normal distribution $\mathcal{N}(0, \sigma^2)$ with

$$(16) \quad \sigma^2 = \frac{1}{\zeta^4} \sum_{b=1}^R \sum_{c=1}^S \sum_{a=1}^K p_{bca} [v(2d_{a|b}^* - \text{Gvar } \mathbf{p}_{(b)}) - \zeta(2d_{a|bc}^* - \text{Gvar } \mathbf{p}_{(bc)})]^2$$

assuming that $\sigma^2 \neq 0$. Here we denote

$$(17) \quad v = \sum_{r=1}^R \sum_{s=1}^S p_{rs+} \text{Gvar } \mathbf{p}_{(rs)}, \quad \zeta = \sum_{r=1}^R p_{r++} \text{Gvar } \mathbf{p}_{(r)},$$

$$\mathbf{p}_{(rs)} = \begin{pmatrix} p_{rs1} \\ p_{rs+} \end{pmatrix},$$

$$\mathbf{p}_{(r)} = \begin{pmatrix} p_{r+1} \\ p_{r++} \end{pmatrix}.$$

$$d_{a|bc}^* = \sum_{k=1}^K \frac{p_{bck}}{p_{bc+}} d_{ka}, \quad d_{a|b}^* = \sum_{k=1}^K \frac{p_{b+k}}{p_{b++}} d_{ka}$$

and f_{111}, \dots, f_{RSK} are observed relative frequencies. Further it holds that

$$\frac{\sqrt{(n)} (\delta_{A|C(B)} - \delta_{A|C(B)})}{\sigma(f_{111}, \dots, f_{RSK})}$$

has the normal distribution $\mathcal{N}(0, 1)$. The asymptotic variance σ^2 is equal to zero if and only if all distributions $\mathbf{p}_{(bc)}$ (with regard to c) belong to the same class of equivalence $\mathbf{p}_{(b)}$ generated by $D(\cdot, \cdot)$, and it holds for all b , i.e. if and only if the coefficient $\delta_{A|C(B)}$ is equal to zero.

The proof is again analogical to that given in Theorem 1 and therefore it will be omitted. It is given in [18].

4. TESTS OF GOODNESS-OF-FIT AND HOMOGENEITY FOR DISTRIBUTIONS OF A GENERALIZED CATEGORICAL VARIABLE

The hypotheses testing for D -variables in standard situations of homogeneity and goodness-of-fit has the high relevance. The nature of the problem represented by the particular matrix \mathbf{D} determines a special interest in types of heterogeneity and departure from the full coincidence of the sampling distribution with the hypothetical assumption. The distance (3) gives a natural measure that considers the differences among distributions in a desired way. We begin by reminding a well known result (see e.g. [8]).

Lemma 1. If $\mathbf{X} = (X_1, \dots, X_K)'$ has the multivariate normal distribution $\mathcal{N}(\mathbf{0}, \Sigma)$, $R(\Sigma) = K$ and $Y = \mathbf{X}'\mathbf{A}\mathbf{X}$ for some symmetric nonnegative definite matrix \mathbf{A} ,

then $\mathcal{L}[Y] = \mathcal{L}[\sum_{i=1}^K \lambda_i Z_i^2]$ where Z_1^2, \dots, Z_K^2 are independent chi-square variables with one degree of freedom each and $\lambda_1, \dots, \lambda_K$ are the eigenvalues of $\Sigma \mathbf{A}$.

Remark. Values of the distribution function $P(Y \leq t)$ can be computed according to Algorithm AS 106 (see [20]).

Theorem 5. (Goodness-of-fit.) If $\mathbf{X} = (X_1, \dots, X_K)'$ has the multinomial distribution $\mathcal{M}_K(n, \mathbf{p})$, $\mathbf{f} = (f_1, \dots, f_K)'$, $f_i = X_i/n$ and $\mathbf{g} \in \mathbf{Q}_K$ is a given vector with nonzero elements, then under the hypothesis $\mathbf{p} = \mathbf{g}$ it holds that

$$(18) \quad \mathcal{L}[nD^2(\mathbf{f}, \mathbf{g})] \rightarrow \mathcal{L}[\sum_{i=1}^{K-1} \lambda_i Z_i^2],$$

where Z_i^2 are independent χ_1^2 variates and λ_i are the eigenvalues of $(\mathbf{D}_{\mathbf{g}^*} - \mathbf{g}^*(\mathbf{g}^*)') \mathbf{D}^*$, $\mathbf{g}^* = (g_1, \dots, g_{K-1})'$, $\mathbf{D}_{\mathbf{g}^*} = \text{diag}\{g_1, \dots, g_{K-1}\}$ and \mathbf{D}^* is given by (4).

Proof. Denote $\mathbf{V}_n = \sqrt{(n)}(\mathbf{f}^* - \mathbf{g}^*)$. It is known that $\mathcal{L}[\mathbf{V}_n] \rightarrow \mathcal{L}[\mathbf{V}]$, where \mathbf{V} has the multivariate normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{D}_{\mathbf{g}^*} - \mathbf{g}^*(\mathbf{g}^*)')$. If we take into account that $R(\mathbf{D}_{\mathbf{g}^*} - \mathbf{g}^*(\mathbf{g}^*)') = K - 1$ and $nD^2(\mathbf{f}, \mathbf{g}) = \mathbf{V}_n' \mathbf{D}^* \mathbf{V}_n$ is a continuous function of \mathbf{V}_n and if we apply Lemma 1, then

$$\mathcal{L}[nD^2(\mathbf{f}, \mathbf{g})] \rightarrow \mathcal{L}[\mathbf{V}' \mathbf{D}^* \mathbf{V}] = \mathcal{L}[\sum_{i=1}^{K-1} \lambda_i Z_i^2]$$

with those λ_i as stated in Theorem 5. □

Lemma 2. Let $\mathbf{X} = (X_1, \dots, X_K)'$ and $\mathbf{Y} = (Y_1, \dots, Y_K)'$ be independent random variables having the multinomial distributions $\mathcal{M}_K(n, \mathbf{p})$ and $\mathcal{M}_K(m, \mathbf{p})$, $\mathbf{p} = (p_1, \dots, p_K)'$; $p_k > 0$ for all k . Denote $\mathbf{f} = n^{-1} \mathbf{X} = (f_1, \dots, f_K)'$, $\mathbf{g} = m^{-1} \mathbf{Y} = (g_1, \dots, g_K)'$, $\mathbf{f}^* = (f_1, \dots, f_{K-1})'$, $\mathbf{g}^* = (g_1, \dots, g_{K-1})'$, $\mathbf{p}^* = (p_1, \dots, p_{K-1})'$. Then

$$(19) \quad \mathcal{L}[\sqrt{(n+m)}(\mathbf{f}^* - \mathbf{g}^*)] \rightarrow \mathcal{N}(\mathbf{0}, [\omega(1-\omega)]^{-1} (\mathbf{D}_{\mathbf{p}^*} - \mathbf{p}^*(\mathbf{p}^*)')),$$

where n and m tends to infinity so that $\lim m/(m+n) = \omega$, $0 < \omega < 1$.

The proof is obvious. □

Theorem 6. (Two-sample problem.) Under the assumptions of Lemma 2 it holds that

$$(20) \quad \mathcal{L}[(n+m)D^2(\mathbf{f}, \mathbf{g})] \rightarrow \mathcal{L}[\sum_{i=1}^{K-1} \lambda_i Z_i^2],$$

where Z_i^2 are independent χ_1^2 variates and λ_i are the eigenvalues of $\omega[(1-\omega)]^{-1} (\mathbf{D}_{\mathbf{p}^*} - \mathbf{p}^*(\mathbf{p}^*)') \mathbf{D}^*$, $\mathbf{D}_{\mathbf{p}^*} = \text{diag}\{p_1, \dots, p_{K-1}\}$, \mathbf{D}^* is given by (4).

Proof. $(n+m)D^2(\mathbf{f}, \mathbf{g}) = (n+m)(\mathbf{f}^* - \mathbf{g}^*)' \mathbf{D}^*(\mathbf{f}^* - \mathbf{g}^*) = \mathbf{H}'_{n+m} \mathbf{D}^* \mathbf{H}_{n+m}$ say. According to Lemma 2

$$\mathcal{L}[\mathbf{H}_{n+m}] \rightarrow \mathcal{L}[\mathbf{H}] = \mathcal{N}(\mathbf{0}, [\omega(1-\omega)]^{-1} (\mathbf{D}_{\mathbf{p}^*} - \mathbf{p}^*(\mathbf{p}^*)')).$$

Further $\mathcal{L}[\mathbf{H}'_{n+m}\mathbf{D}^*\mathbf{H}_{n+m}] \rightarrow \mathcal{L}[\mathbf{H}'\mathbf{D}^*\mathbf{H}]$, because $\mathbf{H}'_{n+m}\mathbf{D}^*\mathbf{H}_{n+m}$ is a continuous function of \mathbf{H}_{n+m} . From Lemma 1 it follows that $\mathcal{L}[\mathbf{H}'\mathbf{D}^*\mathbf{H}] = \mathcal{L}[\sum_{i=1}^{K-1} \lambda_i Z_i^2]$ with those λ_i as stated in Theorem 6. \square

Remark. If we do not know \mathbf{p}^* we may use its estimate $\hat{\mathbf{p}}^*$ with components $\hat{p}_i = (n_i + m_i)/(n + m)$. Suppose that $n_i + m_i = nf_i + mg_i > 0$ ($i = 1, \dots, K$). First we note that $\Sigma = \mathbf{D}_{\mathbf{p}^*} - \mathbf{p}^*(\mathbf{p}^*)'$ and $\hat{\Sigma} = \mathbf{D}_{\hat{\mathbf{p}}^*} - \hat{\mathbf{p}}^*(\hat{\mathbf{p}}^*)'$ are nonsingular matrices and that $\hat{\Sigma} \xrightarrow{P} \Sigma$. Therefore $\hat{\Sigma}^{-1/2} \xrightarrow{P} \Sigma^{-1/2}$. Denote $\mathbf{V}_{n+m} = \sqrt{(\omega(1 - \omega))} \cdot \hat{\Sigma}^{-1/2} \mathbf{H}_{n+m}$. Then $\mathcal{L}[\mathbf{V}_{n+m}] \rightarrow \mathcal{L}[\mathbf{V}] = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mathcal{L}[\mathbf{V}'_{n+m}\mathbf{D}^*\mathbf{V}_{n+m}] \rightarrow \mathcal{L}[\mathbf{V}'\mathbf{D}^*\mathbf{V}] = \mathcal{L}[\sum_{i=1}^{K-1} \lambda_i Z_i^2]$, where λ_i are the eigenvalues of \mathbf{D}^* .

Lemma 3. Let $\mathbf{X}_1, \dots, \mathbf{X}_R, \mathbf{X}_r = (X_{r1}, \dots, X_{rK})'$, $r = 1, \dots, R$, be independent random variables having the distributions $\mathcal{M}_K(n_r, \mathbf{p})$, $r = 1, \dots, R$, $\mathbf{p} = (p_1, \dots, p_K)'$, $p_k > 0$ for all k . Denote $\mathbf{f}_{(r)} = n_r^{-1} \mathbf{X}_r = (f_{1|r}, \dots, f_{K|r})'$, $\mathbf{f}_{(r)}^* = (f_{1|r}^*, \dots, f_{K-1|r}^*)'$, $\mathbf{h} = \sum_{r=1}^R w_r \mathbf{f}_{(r)}$, $\mathbf{h}^* = \sum_{r=1}^R w_r \mathbf{f}_{(r)}^*$, $w_r > 0$ ($r = 1, \dots, R$), $\sum_{r=1}^R w_r = 1$. Let $n_r \rightarrow \infty$ in such a way that $n_r/n \rightarrow \omega_r$, where $n = \sum_{r=1}^R n_r$, $0 < \omega_r < 1$ ($r = 1, \dots, R$), $\sum_{r=1}^R \omega_r = 1$. Then

$$(21) \quad \mathcal{L}[\sqrt{(n)}(\mathbf{f}_{(1)}^* - \mathbf{h}^*, \dots, \mathbf{f}_{(R-1)}^* - \mathbf{h}^*)] \rightarrow \mathcal{N}(\mathbf{0}, \dots, \mathbf{0}),$$

$$\mathbf{B}_{R-1} \otimes \Sigma_{K-1},$$

where

$$(22) \quad \Sigma_{K-1} = \mathbf{D}_{\mathbf{p}^*} - \mathbf{p}^*(\mathbf{p}^*)', \quad \mathbf{D}_{\mathbf{p}^*} = \text{diag}\{p_1, \dots, p_{K-1}\}$$

and \mathbf{B}_{R-1} is a positive definite matrix with elements

$$(23) \quad b_{rr} = \frac{1}{\omega_r} - \frac{2w_r}{\omega_r} + \sum_{i=1}^R \frac{w_i^2}{\omega_i} \quad (r = 1, \dots, R-1)$$

$$b_{rs} = b_{sr} = -\frac{w_r}{\omega_r} - \frac{w_s}{\omega_s} + \sum_{i=1}^R \frac{w_i^2}{\omega_i} \quad (r \neq s = 1, \dots, R-1).$$

Proof. It is easy to show that $\mathcal{L}[\sqrt{(n)}(\mathbf{f}_{(1)}^* - \mathbf{h}^*, \dots, \mathbf{f}_{(R)}^* - \mathbf{h}^*)]$ converges to a normal distribution with expected value vector $(\mathbf{0}, \dots, \mathbf{0})$ and covariance matrix

$$(24) \quad (\mathbf{I}_{R(K-1)} - \mathbf{A}\mathbf{C}) \Sigma_{R(K-1)} (\mathbf{I}_{R(K-1)} - \mathbf{A}\mathbf{C})',$$

where $\mathbf{I}_{R(K-1)}$ denotes the $R(K-1) \times R(K-1)$ unit matrix, $\mathbf{A} = (\mathbf{I}_{K-1}, \dots, \mathbf{I}_{K-1})'$ is a $R(K-1) \times (K-1)$ matrix, $\mathbf{C} = (w_1 \mathbf{I}_{K-1}, \dots, w_R \mathbf{I}_{K-1})$ is a $(K-1) \times R(K-1)$ matrix, \mathbf{I}_{K-1} denotes the $(K-1) \times (K-1)$ unit matrix and

$$\Sigma_{R(K-1)} = \left\| \begin{array}{cccc} \omega_1^{-1} \Sigma_{K-1}, & \mathbf{0}, & \dots, & \mathbf{0} \\ \mathbf{0}, & \omega_2^{-1} \Sigma_{K-1}, & \dots, & \mathbf{0} \\ \dots, & \dots, & \dots, & \dots \\ \mathbf{0}, & \mathbf{0}, & \dots, & \omega_R^{-1} \Sigma_{K-1} \end{array} \right\|$$

is a $R(K-1) \times R(K-1)$ matrix with Σ_{K-1} given by (22).

After some algebra we find that (24) is equal to $\mathbf{B}_R \otimes \Sigma_{K-1}$, where \mathbf{B}_R has the elements given by (23) for $r, s = 1, \dots, R$ and the symbol “ \otimes ” stands for the Kronecker product of matrices. And now it is easy to conclude that the statement (21) holds.

The matrix \mathbf{B}_{R-1} arises from \mathbf{B}_R by omitting the last row and column. \mathbf{B}_{R-1} is positive semidefinite because $\mathbf{B}_{R-1} \otimes \Sigma_{K-1}$ is a covariance matrix. \mathbf{B}_{R-1} is also nonsingular. It follows from the fact that $\mathbf{B}_R = \mathbf{U}_R \mathbf{U}_R'$.

$$\mathbf{U}_R = \left\| \begin{array}{cccc} (1-w_1)/\sqrt{\omega_1}, & -w_2/\sqrt{\omega_2}, & \dots, & -w_R/\sqrt{\omega_R} \\ -w_1/\sqrt{\omega_1}, & -w_2/\sqrt{\omega_2}, & \dots, & (1-w_R)/\sqrt{\omega_R} \end{array} \right\|,$$

$$R(\mathbf{B}_R) = R(\mathbf{U}_R \mathbf{U}_R') = R(\mathbf{U}_R) = R - 1, \quad \text{since } |\mathbf{U}_R| = 0,$$

$$|\mathbf{U}_{R-1}| = w_R \prod_{r=1}^{R-1} (1/\sqrt{\omega_r}) > 0.$$

The last row (column) of \mathbf{B}_R is a linear combination of the preceding rows (columns). If we omit them the rank will not change. Hence \mathbf{B}_{R-1} is positive definite. \square

Theorem 7. (*R-sample problem, one-way ANOVA for D-variables.*) Under the assumptions of Lemma 3 it holds that

$$(25) \quad \mathcal{L}\left[n \sum_{r=1}^R w_r D^2(\mathbf{f}_{(r)}, \mathbf{h})\right] \rightarrow \mathcal{L}\left[\sum_{i=1}^{(R-1)(K-1)} \lambda_i Z_i^2\right],$$

where Z_i^2 are independent χ_1^2 variates and λ_i are the eigenvalues of the matrix $(\mathbf{B}_{R-1} \otimes \Sigma_{K-1})(\mathbf{W}_{R-1} \otimes \mathbf{D}^*)$, where $\mathbf{B}_{R-1}, \Sigma_{K-1}, \mathbf{W}_{R-1}$ are positive definite matrices given by (23), (22) and

$$(26) \quad W_{rr} = w_r(w_R + w_r)/w_R \quad (r = 1, \dots, R-1)$$

$$W_{rs} = W_{sr} = w_r w_s / w_R \quad (r \neq s)$$

\mathbf{D}^* is a nonnegative definite matrix with elements given by (4).

Proof. After some algebraic manipulation we get that

$$\begin{aligned} & n \sum_{r=1}^R w_r D^2(\mathbf{f}_{(r)}, \mathbf{h}) = \\ & = n(\mathbf{f}_{(1)}^* - \mathbf{h}^*, \dots, \mathbf{f}_{(R-1)}^* - \mathbf{h}^*)' (\mathbf{W}_{R-1} \otimes \mathbf{D}^*) (\mathbf{f}_{(1)}^* - \mathbf{h}^*, \dots, \mathbf{f}_{(R-1)}^* - \mathbf{h}^*) = \\ & = \mathbf{H}'_n \mathbf{M} \mathbf{H}_n \quad (\text{say}). \end{aligned}$$

The matrix \mathbf{W}_{R-1} is positive definite because its principal minors are equal to $w_1 \dots w_i (w_1 + \dots + w_i + w_R)/w_R$ for $i = 1, \dots, R-1$. The matrix \mathbf{M} is positive definite (semidefinite) if \mathbf{D}^* is positive definite (semidefinite). According to Lemma 3 $\mathcal{L}[\mathbf{H}_n] \rightarrow \mathcal{L}[\mathbf{H}] = \mathcal{N}((\mathbf{0}, \dots, \mathbf{0}), \mathbf{B}_{R-1} \otimes \Sigma_{K-1})$ and according to Lemma 1 $\mathcal{L}[\mathbf{H}'\mathbf{M}\mathbf{H}] = \mathcal{L}\left[\sum_{i=1}^{(R-1)(K-1)} \lambda_i Z_i^2\right]$, where λ_i are the eigenvalues of $(\mathbf{B}_{R-1} \otimes \Sigma_{K-1}) \mathbf{M}$.

Remark. If the vector \mathbf{p}^* is unknown we may use its estimate $\hat{\mathbf{p}}^*$ with components $\hat{p}_k = \sum_{r=1}^R n_{rk}/n$, where n_{rk} is an observed value of X_{rk} . Suppose $\hat{p}_k > 0$ for $k = 1, \dots, K$. Σ_{K-1} and its estimate $\hat{\Sigma}_{K-1}$ are nonsingular, $\hat{\Sigma}_{K-1} \xrightarrow{P} \Sigma_{K-1}$, hence $\hat{\Sigma}_{K-1}^{-1/2} \xrightarrow{P} \Sigma_{K-1}^{-1/2}$ and

$$\begin{aligned} \mathcal{L}[(\mathbf{B}_{R-1} \otimes \Sigma_{K-1})^{-1/2} \mathbf{H}_n] &= \mathcal{L}[\mathbf{V}_n] \rightarrow \mathcal{L}[(\mathbf{B}_{R-1} \otimes \Sigma_{K-1})^{-1/2} \mathbf{H}] = \\ &= \mathcal{L}[\mathbf{V}] = \mathcal{N}(\mathbf{0}, \dots, \mathbf{0}), \mathbf{1}_{(R-1)(K-1)} \end{aligned}$$

and $\mathcal{L}[\mathbf{V}'_n \mathbf{M} \mathbf{V}_n] \rightarrow \mathcal{L}[\mathbf{V}' \mathbf{M} \mathbf{V}] = \mathcal{L}[\sum_{i=1}^{(R-1)(K-1)} \lambda_i \mathbf{Z}_i^2]$, where λ_i are the eigenvalues of the matrix \mathbf{M} .

Theorem 6 is a special case of Theorem 7. However, it has been separated for its simpler formulation and computational aspects as well as for the frequent occurrence of the two-sample case in practice. The two-sample statistic can also be used in R -sample problem in combination with a simultaneous testing or in the case of rejecting homogeneity hypothesis if we want to investigate the distributions pair-wise to find out which populations differ.

The investigation of residuals after the hypothesis of homogeneity ($R \geq 2$) is rejected is an important part of the data analysis. For the purposes of the next theorem let us denote

$$\begin{aligned} \mathbf{f}^{(r)} &= (n - n_{r+})^{-1} \sum_{u \neq r} n_{u+} \mathbf{f}_{(u)} = (f_1^{(r)}, \dots, f_K^{(r)})', \\ \mathbf{p}^{(r)} &= (n - n_{r+})^{-1} \sum_{u \neq r} n_{u+} \mathbf{p}_{(u)} = (p_1^{(r)}, \dots, p_K^{(r)})', \\ \mathbf{q}_{(r)} &= \mathbf{D}(\mathbf{f}_{(r)} - \mathbf{f}^{(r)}) = (q_{1|r}, \dots, q_{K|r})' \end{aligned}$$

($\mathbf{q}_{(r)}$ is a residual characteristic of the contingency table row).

Theorem 8. (Residual analysis of the departure from homogeneity.) Let $\mathbf{f}_{(1)}, \dots, \mathbf{f}_{(R)}$ are distributed independently as $\mathcal{L}[n_{r+}, \mathbf{f}_{(r)}] = \mathcal{M}(n_{r+}, \mathbf{p}_{(r)})$, $p_{k|r} > 0$ for all k and r . If $n_{r+} \rightarrow \infty$ in such a way that $n_{r+}/n \rightarrow \omega_r$, $0 < \omega_r < 1$, $r = 1, \dots, R$, $n = \sum n_{r+}$, then it holds that

$$\begin{aligned} \mathcal{L}[\sqrt{(n)} (\mathbf{q}_{(r)} - \mathbf{D}(\mathbf{p}_{(r)} - \mathbf{p}^{(r)}), \mathbf{q}_{(s)} - \mathbf{D}(\mathbf{p}_{(s)} - \mathbf{p}^{(s)}))] &\rightarrow \\ &\rightarrow \mathcal{N}(\mathbf{0}, \mathbf{0}), \mathbf{A} \Sigma \mathbf{A}' \end{aligned}$$

where

$$\begin{aligned} \mathbf{A} &= \begin{Bmatrix} \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{Bmatrix}, \quad \mathbf{A} \Sigma \mathbf{A}' = \begin{Bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_2 & \mathbf{A}_3 \end{Bmatrix}, \\ (27) \quad (\mathbf{A}_1)_{ab} &= \omega_r^{-1} \left\{ \sum_{k=1}^K p_{k|r} (d_{ak} - d_{a|r}^*) (d_{bk} - d_{b|r}^*) \right\} + \\ &+ (1 - \omega_r)^{-2} \left\{ \sum_{u \neq r} \omega_u \left[\sum_{k=1}^K p_{k|u} (d_{ak} - d_{a|u}^*) (d_{bk} - d_{b|u}^*) \right] \right\}, \end{aligned}$$

$$\begin{aligned}
(28) \quad (A_2)_{ab} &= [(1 - \omega_r)(1 - \omega_s)]^{-1} \left\{ \sum_{u \neq r, s} \omega_u \left[\sum_{k=1}^K p_{k|u} (d_{ak} - d_{a|u}^*) (d_{bk} - d_{b|u}^*) \right] \right\} - \\
&\quad - (1 - \omega_s)^{-1} \left\{ \sum_{k=1}^K p_{k|r} (d_{ak} - d_{a|r}^*) (d_{bk} - d_{b|r}^*) \right\} - \\
&\quad - (1 - \omega_r)^{-1} \left\{ \sum_{k=1}^K p_{k|s} (d_{ak} - d_{a|s}^*) (d_{bk} - d_{b|s}^*) \right\}. \\
(29) \quad (A_3)_{ab} &= \omega_s^{-1} \left\{ \sum_{k=1}^K p_{k|s} (d_{ak} - d_{a|s}^*) (d_{bk} - d_{b|s}^*) \right\} + \\
&\quad + (1 - \omega_s)^{-2} \left\{ \sum_{u \neq s} \omega_u \left[\sum_{k=1}^K p_{k|u} (d_{ak} - d_{a|u}^*) (d_{bk} - d_{b|u}^*) \right] \right\}, \\
&\quad a, b = 1, \dots, K \text{ and } d_{a|r}^* = \sum_j d_{aj} p_{j|r}.
\end{aligned}$$

Proof. The asymptotic normality follows from the properties of frequencies. In the first step we can see that

$$\mathcal{L}[\sqrt{(n)}(\mathbf{f}_{(r)} - \mathbf{f}^{(r)} - \mathbf{p}_{(r)} + \mathbf{p}^{(r)}, \mathbf{f}_{(s)} - \mathbf{f}^{(s)} - \mathbf{p}_{(s)} + \mathbf{p}^{(s)})'] \rightarrow \mathcal{N}((\mathbf{0}, \mathbf{0})', \mathbf{\Sigma}),$$

where $\mathbf{\Sigma}$ is a block matrix

$$\begin{aligned}
\mathbf{\Sigma} &= \begin{bmatrix} \mathbf{\Sigma}_1 & \mathbf{\Sigma}_2 \\ \mathbf{\Sigma}_2 & \mathbf{\Sigma}_3 \end{bmatrix}, \\
(\Sigma_1)_{aa} &= \omega_r^{-1} p_{a|r} (1 - p_{a|r}) + (1 - \omega_r)^{-2} \sum_{u \neq r} \omega_u p_{a|u} (1 - p_{a|u}), \\
(\Sigma_1)_{ab} &= -\omega_r^{-1} p_{a|r} p_{b|r} - (1 - \omega_r)^{-2} \sum_{u \neq r} \omega_u p_{a|u} p_{b|u}, \\
(\Sigma_2)_{aa} &= [(1 - \omega_r)(1 - \omega_s)]^{-1} \sum_{u \neq r, s} \omega_u p_{a|u} (1 - p_{a|u}) - \\
&\quad - (1 - \omega_s)^{-1} p_{a|r} (1 - p_{a|r}) - (1 - \omega_r)^{-1} p_{a|s} (1 - p_{a|s}), \\
(\Sigma_2)_{ab} &= -[(1 - \omega_r)(1 - \omega_s)]^{-1} \sum_{u \neq r, s} \omega_u p_{a|u} p_{b|u} + \\
&\quad + (1 - \omega_s)^{-1} p_{a|r} p_{b|r} + (1 - \omega_r)^{-1} p_{a|s} p_{b|s}.
\end{aligned}$$

The matrix $\mathbf{\Sigma}_3$ is the same as $\mathbf{\Sigma}_1$, only we must substitute r by s . The rest follows from the fact that

$$\begin{aligned}
\mathbf{A}(\mathbf{f}_{(r)} - \mathbf{f}^{(r)} - \mathbf{p}_{(r)} + \mathbf{p}^{(r)}, \mathbf{f}_{(s)} - \mathbf{f}^{(s)} - \mathbf{p}_{(s)} + \mathbf{p}^{(s)})' &= \\
&= (\mathbf{q}_{(r)} - \mathbf{D}(\mathbf{p}_{(r)} - \mathbf{p}^{(r)}), \mathbf{q}_{(s)} - \mathbf{D}(\mathbf{p}_{(s)} - \mathbf{p}^{(s)}))'. \quad \square
\end{aligned}$$

The next corollary gives the tool for examining the significance of residuals from hypothesis of homogeneity.

Corollary 1. Under the assumptions of Theorem 8 and under $\mathbf{H}_0: \mathbf{p}_{(r)} = \mathbf{p}$, $r = 1, \dots, R$, it holds that

$$\mathcal{L}[\sqrt{(n)} q_{k|r}] \rightarrow \mathcal{N}(0, \sigma_q^2),$$

where

$$\sigma_q^2 = [\omega_r(1 - \omega_r)]^{-1} \sum_{j=1}^K p_j (d_{jk} - \sum_{a=1}^K p_a d_{ak})^2$$

and

$$\mathcal{L}[\sqrt{(n)} q_{k|r}/\hat{\sigma}_q] \rightarrow \mathcal{N}(0, 1),$$

where

$$\hat{\sigma}_q^2 = \{n^2/[n_{r+}(n - n_{r+})]\} \sum_{j=1}^K f_{+j} d_{jk} - \sum_{a=1}^K f_{+a} d_{ak}^2$$

and $f_{+j} = n^{-1} \sum_{j|s} n_{s+}$ is supposed to be positive for all j .

The statistics $\sqrt{(n)} q_{k|r}/\hat{\sigma}_q$ may be used for the testing \mathbf{H}_0 against $(\mathbf{p}_{(r)} - \mathbf{p}^{(r)})' \mathbf{D}_k \neq 0$, $\mathbf{D}_k = (d_{1k}, \dots, d_{Kk})'$. Instead of it we may use the statistics $\sqrt{(n)} t_{k|r}/\hat{\sigma}_r$, where

$$t_{k|r} = \sum_{j=1}^K (f_{j|r} - f_{+j}) d_{jk} = (1 - f_{r+}) q_{k|r},$$

$$\hat{\sigma}_r = (1 - f_{r+}) \hat{\sigma}_q.$$

Notice that σ_q^2 and $\hat{\sigma}_q^2$ are equal to zero if and only if $\mathbf{D}_k = \mathbf{0}$. Of course such a column (and the corresponding row) of the matrix \mathbf{D} would be excluded in advance, so that we can see that σ_q^2 and $\hat{\sigma}_q^2$ are always greater than zero.

Corollary 2. Under the assumptions of Theorem 8

$$\mathcal{L}[\sqrt{(n)} (q_{k|r} - q_{l|s} - \mathbf{E}(q_{k|r} - q_{l|s}))] \rightarrow \mathcal{N}(0, \sigma^2)$$

provided that $\sigma^2 \neq 0$, where σ^2 is given in the following way:

1. $k = l, \quad r \neq s \Rightarrow \sigma^2 = (A_1)_{kk} + (A_3)_{kk} - 2(A_2)_{kk},$
2. $k \neq l, \quad r = s \Rightarrow \sigma^2 = (A_1)_{kk} + (A_1)_{ll} - 2(A_1)_{kl},$
3. $k \neq l, \quad r \neq s \Rightarrow \sigma^2 = (A_1)_{kk} + (A_3)_{ll} - 2(A_2)_{kl}.$

$\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3$ are described in (27)–(29) and

$$\mathbf{E}(q_{k|r} - q_{l|s}) = d_{k|r}^* - (n - n_{r+})^{-1} \sum_{u \neq r} n_{u+} d_{k|u}^* - d_{l|s}^* + (n - n_{s+})^{-1} \sum_{v \neq s} n_{v+} d_{l|v}^*.$$

The last part of this article deals with marginal homogeneity in square tables.

Lemma 4. Let $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_K; \mathbf{D}\}$, $\mathbf{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_K; \mathbf{D}\}$ be D -variables with the same matrix \mathbf{D} . The square table $\mathbf{A} \times \mathbf{B}$ arises as a sample from the multinomial distribution $\mathcal{M}(n; p_{11}, \dots, p_{KK})$, $p_{ij} > 0$ for all i, j . Denote $\mathbf{p}_\mathbf{A} = (p_{1+}, \dots, p_{K+})'$, $\mathbf{p}_\mathbf{B} = (p_{+1}, \dots, p_{+K})'$. Let n_{ij} be the observed frequency in the cell (i, j) and denote $\mathbf{f}_\mathbf{A} = (f_{1+}, \dots, f_{K+})'$, $\mathbf{f}_\mathbf{B} = (f_{+1}, \dots, f_{+K})'$, $\mathbf{f}_\mathbf{A}^* = (f_{1+}, \dots, f_{K-1,+})'$, $\mathbf{f}_\mathbf{B}^* = (f_{+1}, \dots, f_{+,K-1})'$, $f_{i+} = n_{i+}/n$, $f_{+i} = n_{+i}/n$. Under the condition $\mathbf{p}_\mathbf{A} = \mathbf{p}_\mathbf{B}$ it holds that

$$(30) \quad \mathcal{L}[\sqrt{(n)} (\mathbf{f}_\mathbf{A}^* - \mathbf{f}_\mathbf{B}^*)] \rightarrow \mathcal{N}(\mathbf{0}, \Sigma),$$

where Σ is a $(K - 1) \times (K - 1)$ positive definite matrix with elements

$$(31) \quad \begin{aligned} \sigma_{ii} &= p_{i+} + p_{+i} - 2p_{ii} \quad (i = 1, \dots, K - 1) \\ \sigma_{ij} &= -(p_{ij} + p_{ji}) \quad (i \neq j). \end{aligned}$$

For a proof of this result, see e.g. [1], p. 219.

Theorem 9. (*Marginal homogeneity in square tables.*) Under the assumption of Lemma 4 the asymptotic distribution of $D^2(\mathbf{f}_A, \mathbf{f}_B)$ is given by

$$(32) \quad \mathcal{L}[nD^2(\mathbf{f}_A, \mathbf{f}_B)] \rightarrow \mathcal{L}\left[\sum_{i=1}^{K-1} \lambda_i Z_i^2\right],$$

where Z_i^2 are independent χ_1^2 variables and λ_i are the eigenvalues of the matrix $\Sigma \mathbf{D}^*$, where Σ is defined as in Lemma 4 and \mathbf{D}^* is given by (4).

Proof. The desired result immediately follows from Lemma 1 and 4 in the same way as in the proof of Theorem 6. \square

Remark. If we do not know the probabilities p_{ij} we may use $\hat{p}_{ij} = n_{ij}/n$. Under the conditions of Lemma 4 and assuming that $\hat{p}_{ij} > 0$ ($i, j = 1, \dots, K$) it holds that $\mathcal{L}[\mathbf{V}_n] \rightarrow \mathcal{L}[\mathbf{V}] = \mathcal{N}(\mathbf{0}, \mathbf{I})$, where $\mathbf{V}_n = \sqrt{(n)} \hat{\Sigma}^{-1/2}(\mathbf{f}_A^* - \mathbf{f}_B^*)$ and $\hat{\Sigma}$ is an estimate of Σ with \hat{p}_{ij} instead of p_{ij} . Then $\mathcal{L}[\mathbf{V}_n \mathbf{D}^* \mathbf{V}_n] \rightarrow \mathcal{L}[\mathbf{V} \mathbf{D}^* \mathbf{V}] = \mathcal{L}\left[\sum_{i=1}^{K-1} \lambda_i Z_i^2\right]$, where λ_i are the eigenvalues of \mathbf{D}^* .

CONCLUSIONS

The statistical theory of asymptotic behaviour of measures derived from the D -model provides a *unified approach* to different types of situations in which the researcher wants to specify the relations among the categories of variables. Important special cases as nominal, ordinal, cardinal or geographical variables are also included. Also analyses of special experimental or decisional data can be done this way if the elements of the matrix \mathbf{D} are interpreted as the loss function values.

(Received October 27, 1986.)

REFERENCES

- [1] J. Anděl: Matematická statistika (Mathematical Statistics). SNTL, Praha 1978.
- [2] L. A. Goodman and W. H. Kruskal: Measures of association for cross classification IV: simplification of asymptotic variances. J. Amer. Statist. Assoc. *67* (1972), 415–421.
- [3] C. R. Rao: Lineární metody statistické indukce a jejich aplikace (Linear Statistical Inference and Its Applications). Academia, Praha 1978.
- [4] C. R. Rao: Diversity: its measurement, decomposition, apportionment and analysis. *Sankhyā* *44*, Series A (1982), 1–22.
- [5] C. R. Rao: Gini-Simpson index of diversity: a characterization, generalization and applications. *Utilitas Mathematica* *21* B (1982), 273–282.
- [6] C. R. Rao: Analysis of diversity: a unified approach. *Statistical Decision Theory and Related Topics III*, 2 (1982), 233–250.
- [7] C. R. Rao: Diversity and dissimilarity coefficients: a unified approach. *Theoretical Population Biology* *21* (1982), 24–43.
- [8] H. Ruben: Probability content of regions under spherical normal distributions, IV: the distribution of homogeneous and nonhomogeneous quadratic functions of normal variables. *Ann. Math. Statist.* *33* (1962), 542–570.

- [9] J. Řehák: Základní deskriptivní míry pro rozložení ordinálních dat (Basic descriptive measures for the distribution of ordinal data). Sociologický časopis 12 (1976), 4, 416—431.
- [10] J. Řehák and I. Loučková: Klasické mnohorozměrné škálování (Aplikace metody DISTAN) (Classical multidimensional scaling (application of the DISTAN method)). Sociologický časopis 19 (1983), 535—554.
- [11] J. Řehák and I. Loučková: Komparace podmíněných distribucí kontingenční tabulky klasickým mnohorozměrným škálováním (Comparison of conditional distributions in a contingency table by classical multidimensional scaling). In: Robust 84 (J. Antoch, J. Jurečková, eds.), JČMF, Praha 1984, 100—104.
- [12] J. Řehák and B. Řeháková: Basic characteristics for finite-valued variables and a distance analysis of their distributions. 9th World Congress of Sociology, Uppsala 1978.
- [13] J. Řehák and B. Řeháková: Distanční přístup k analýze kategorizovaných dat a jeho aplikace na problém shody (Distance approach to analysis of categorical data and its application to the problem of goodness-of-fit). In: Robust 82 (J. Antoch, J. Jurečková, eds.), JČMF, Praha 1982, 76—80.
- [14] J. Řehák and B. Řeháková: Parciální asociační koeficienty v kontingenčních tabulkách (Partial association coefficients in contingency tables). In: Robust 84 (J. Antoch, J. Jurečková, eds.), JČMF, Praha 1984, 105—108.
- [15] J. Řehák and B. Řeháková: Classifications with relations: a model for the description of distributions and their distances. Kybernetika 22 (1986), 2, 158—175.
- [16] J. Řehák and B. Řeháková: Vícenásobná a parciální asociace v kontingenčních tabulkách (Multiple and partial association in contingency tables). Sociologický časopis 22 (1986), 374—394.
- [17] B. Řeháková: Koeficienty parciální asociace pro obecný typ kategorizované proměnné (Partial Association Coefficients for General Type of Categorical Variable). Unpublished Ph. D. Thesis 1982.
- [18] B. Řeháková: Model a metoda pro analýzu kategorizovaných dat s relacemi (A Model and Method for the Analysis of Categorical Data with Relations). Ph. D. Dissertation. Department of Probability and Mathematical Statistics, Charles University, Praha 1985.
- [19] B. Řeháková: Asymptotické testy hypotéz pro rozdělení D-proměnné (Asymptotic tests of hypotheses for distributions of D-variable). In: Robust 86 (J. Antoch, J. Jurečková, eds.), JČMF, Praha 1986, 121—124.
- [20] J. Sheil and I. O'Muircheartaigh: The distribution of non-negative quadratic forms in normal variables. Applied Statistics 26 (1977), 92—98.

RNDr. Jan Řehák, RNDr. Blanka Řeháková, CSc., Geografický ústav ČSAV (Geographical Institute — Czechoslovak Academy of Sciences), Wenzigova 7, 120 00 Praha 2, Czechoslovakia.