

Analysis of Variance based on the Hamming Distance

BY HILDETE PRISCO PINHEIRO *

Department of Statistics

State University of Campinas

FRANÇOISE SEILLIER-MOISEIWITSCH

Department of Biostatistics

University of North Carolina at Chapel Hill

PRANAB KUMAR SEN

Department of Biostatistics

University of North Carolina at Chapel Hill

Abstract

The interest here is the comparison of sequences within or between groups. Sequences are considered on an individual basis, i.e., all possible pairwise comparisons within and across groups are performed. We develop a categorical analysis-of-variance framework based on Hamming distances, the proportion of positions at which two aligned sequences differ, and estimate the variability between, within and across groups. We assume that the sequences are independent, but the positions may not be. In this context U-statistics are utilized to represent the average distance between and within groups as well as the overall distance. The total sum of squares is decomposed into within-, between- and across-group sums of squares. The latter term is new: it does not appear in the classical set-up. Generalized-U-statistics theory (Puri & Sen, 1971; Lee, 1990; Sen & Singer, 1993) is used to find the asymptotic distributions of each sum of squares. Test statistics are developed to assess homogeneity among groups.

1. Introduction The focus here lies in the comparison of sequences. The sequences are considered on an individual basis in the sense that they are compared to each other: all possible pairwise comparisons within and across groups are performed.

*This research was funded in part by CAPES, FAPESP (Brazilian Institutions), the National Science Foundation, the American Foundation for AIDS Research and the National Institutes of Health.

Key words and phrases: Analysis of Variance, Categorical Data, Hamming Distance, U-statistics

We develop an analysis-of-variance framework based on Hamming distances and estimate the variability between, within and across groups (Section 2). In the within sum of squares, we are estimating the variability among individuals within a group around the average distance within this group. In the across sum of squares, we are estimating the variability of individuals across two groups with respect to the average distance between those groups. In the between sum of squares, we estimate the variability in the group average distances around the overall distance.

Weir (1990a) describes an analysis of variance for the genetic variation in the population, in particular for the amount of observed *heterozygosity*. The variance of the estimate of the average heterozygosity is broken down to show the contribution of populations, loci and individuals by setting out the calculations in a framework similar to that of an analysis of variance. Our situation is a little different because we would like to construct a categorical analysis of variance based on Hamming distances (Seillier-Moisewitsch et al., 1994 and references therein), assuming that the sequences are independent, but the positions may not be. The Hamming distance is the proportion of positions at which two aligned sequences differ.

In this context U-statistics are utilized to represent the average distance between and within groups as well as the overall distance (Sections 3 and 4). The total sum of squares is decomposed into within-, between- and across-group sums of squares. The latter term is new: it does not appear in the classical set-up. Generalized-U-statistics theory (Puri & Sen, 1971; Lee, 1990; Sen & Singer, 1993) is used to find the asymptotic distributions of each sum of squares. In Section 5 test statistics are developed to assess homogeneity among groups. The power of the tests are discussed in Section 6. Finally, a data analysis is shown in Section 7.

2. The Total Sum of Squares and its decomposition Let $\mathbf{X}_i^g = (X_{i1}^g, X_{i2}^g, \dots, X_{ik}^g)'$ be a random vector representing sequence i of group g . Suppose $i = 1, \dots, N$, $k = 1, \dots, K$ and $g = 1, \dots, G$. So, X_{ik}^g represents either the amino acid or the nucleotide present at position k of sequence i in group g (e.g., at the nucleotide level, $x_{ik}^g \in \{\mathbf{A}, \mathbf{C}, \mathbf{T}, \mathbf{G}\}$).

Consider $\mathbf{X}_i^{g_1}$ and $\mathbf{X}_j^{g_2}$.

Definition 1

The *Hamming Distance* $D_{ij}^{(g_1 g_2)}$ is a descriptive statistic for sequence comparison defined by

$$\begin{aligned} D_{ij}^{(g_1, g_2)} &= \frac{1}{K} \sum_{k=1}^K I(X_{ik}^{g_1} \neq X_{jk}^{g_2}) \\ &= \frac{1}{K} \times (\text{number of positions where } X_i^{g_1} \text{ and } X_j^{g_2} \text{ differ}), \end{aligned} \tag{2.1}$$

and when $g_1 = g_2 = g$,

$$D_{ij}^g = \frac{1}{K} \sum_{k=1}^K I(X_{ik}^g \neq X_{jk}^g).$$

■

Let $\theta_k^g = P\{X_{ik}^g \neq X_{jk}^g\}$ and $\bar{\theta}^g = \frac{1}{K} \sum_{k=1}^K \theta_k^g$. Then,

$$E[D_{ij}^g] = \frac{1}{K} \sum_{k=1}^K E[I(X_{ik}^g \neq X_{jk}^g)] = \frac{1}{K} \sum_{k=1}^K \theta_k^g = \bar{\theta}^g.$$

Define the average distance within a group as

$$\bar{D}^g = \binom{N}{2}^{-1} \sum_{1 \leq i < j \leq N} D_{ij}^g = \binom{N}{2}^{-1} \frac{1}{K} \sum_{1 \leq i < j \leq N} \sum_{k=1}^K I(X_{ik}^g \neq X_{jk}^g)$$

which is a U-statistic of degree 2 (Lee, 1990). The average distance between two groups is

$$\bar{D}^{(g_1, g_2)} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N D_{ij}^{(g_1, g_2)} = \frac{1}{N^2 K} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^K I(X_{ik}^{g_1} \neq X_{jk}^{g_2})$$

which is a two-sample U-statistics of degree (1,1) (Hoeffding, 1948; Puri & Sen, 1971; Lee, 1990). The overall distance is

$$\begin{aligned} \bar{D} &= \left[G \binom{N}{2} + N^2 \binom{G}{2} \right]^{-1} \left(\sum_{g=1}^G \sum_{1 \leq i < j \leq N} D_{ij}^g + \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N D_{ij}^{(g_1, g_2)} \right) \\ &= \binom{NG}{2}^{-1} \left(\sum_{g=1}^G \binom{N}{2} \bar{D}^g + \sum_{1 \leq g_1 < g_2 \leq G} N^2 \bar{D}^{(g_1, g_2)} \right) \end{aligned}$$

which is a linear combination of U-statistics.

The Total Sum of Squares can be decomposed as

$$\begin{aligned} TSS &= \sum_{g=1}^G \sum_{1 \leq i < j \leq N} (D_{ij}^g - \bar{D}^g)^2 + \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N (D_{ij}^{(g_1, g_2)} - \bar{D}^{(g_1, g_2)})^2 \quad (2.2) \\ &= \sum_{g=1}^G \sum_{1 \leq i < j \leq N} (D_{ij}^g - \bar{D}^g)^2 + \sum_{g=1}^G \sum_{1 \leq i < j \leq N} (\bar{D}^g - \bar{D})^2 \\ &\quad + \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N (D_{ij}^{(g_1, g_2)} - \bar{D}^{(g_1, g_2)})^2 + \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N (\bar{D}^{(g_1, g_2)} - \bar{D})^2 \\ &= WSS + BSS + AWSS + ABSS \end{aligned}$$

where *WSS* stands for Within Sum of Squares, *BSS* for Between Sum of Squares, *AWSS* for Across Within Sum of Squares and *ABSS* for Across Between Sum of Squares

3. Connections Between Sums of Squares and U-statistics Since we have G groups of N sequences, we can disregard the group clustering and think of the sequences as a random sample of size NG . Then

$$\begin{aligned} TSS &= \sum_{1 \leq i < j \leq NG} (D_{ij} - \bar{D})^2 \\ &= \left(\frac{NG(NG-1)}{2} - 1 \right) \left(\frac{NG(NG-1)}{2} \right)^{-1} \sum_{\substack{i < j, i' < j' \\ i \leq i' \text{ or } j \leq j'}} \frac{(D_{ij} - D_{i'j'})^2}{2} \end{aligned} \quad (3.3)$$

$$\begin{aligned} WSS &= \sum_{g=1}^G \sum_{1 \leq i < j \leq N} (D_{ij}^g - \bar{D}^g)^2 \\ &= \left(\frac{N(N-1)}{2} - 1 \right) \left(\frac{N(N-1)}{2} \right)^{-1} \sum_{g=1}^G \sum_{\substack{i < j, i' < j' \\ i \leq i' \text{ or } j \leq j'}} \frac{(D_{ij}^g - D_{i'j'}^g)^2}{2} \end{aligned} \quad (3.4)$$

and

$$\begin{aligned} AWSS &= \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N (D_{ij}^{(g_1, g_2)} - \bar{D}^{(g_1, g_2)})^2 \\ &= (N^2 - 1) \binom{N^2}{2}^{-1} \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N \sum_{i'=1}^N \sum_{\substack{j'=1 \\ i \leq i' \text{ or } j \leq j'}}^N \frac{(D_{ij}^{(g_1, g_2)} - D_{i'j'}^{(g_1, g_2)})^2}{2} \end{aligned} \quad (3.5)$$

The above sums of squares can also be expressed as linear combinations of U-statistics (Pinheiro, 1997). For instance, WSS is a linear combination of one-sample U-statistics of degrees 3 and 4, and $AWSS$ is a linear combination of two-sample U-statistics of degrees (2,2) and (2,1).

4. Asymptotic Distributions and decompositions of U-statistics

Let U_n be a U-statistic of degree m with kernel $\phi(X_1, \dots, X_m)$ and $E(U_n) = \theta(F) = \theta$.

$$U_n = U(X_1, \dots, X_n) = \binom{n}{m}^{-1} \sum_{1 \leq i_1 < \dots < i_m \leq n} \phi(X_{i_1}, \dots, X_{i_m}), \quad n \geq m \quad (4.6)$$

where

$$\theta(F) = E_F\{\phi(X_1, \dots, X_m)\} = \int \dots \int \phi(x_1, \dots, x_m) dF(x_1) \dots dF(x_m)$$

Let

$$\Psi_c(x_1, \dots, x_c) \equiv E\{\phi(x_1, \dots, x_c, X_{c+1}, \dots, X_m)\} \quad (4.7)$$

$$\psi_c(x_1, \dots, x_c) \equiv \mathbb{E}\{\phi(x_1, \dots, x_c, X_{c+1}, \dots, X_m) - \theta\}, \quad (4.8)$$

$$\xi_c \equiv \mathbb{E}\{\psi_c^2(X_1, \dots, X_c)\} = \mathbb{E}\{\Psi_c^2(X_1, \dots, X_c)\} - \theta^2 \quad \text{and} \quad \xi_0 \equiv 0. \quad (4.9)$$

Theorem 1

The function Ψ_c defined in (4.7) has the properties

- (i) $\Psi_c(x_1, \dots, x_c) = \mathbb{E}\{\Psi_d(x_1, \dots, x_c, X_{c+1}, \dots, X_d)\}$ for $1 \leq c < d \leq m$,
- (ii) $\mathbb{E}\{\Psi_c(x_1, \dots, x_c)\} = \mathbb{E}\{\phi(X_1, \dots, X_m)\}$. ■

The proof appears in Lee (1990, p. 11).

By (4.6) and (4.8),

$$\text{Var}(U_n) = \binom{n}{m}^{-2} \sum_{c=0}^m \sum^{(c)} \text{Cov}\{\phi(X_{i_1}, \dots, X_{i_m})\phi(X_{j_1}, \dots, X_{j_m})\}$$

where $\sum^{(c)}$ stands for summation over all subscripts such that

$$1 \leq i_1 < i_2 < \dots < i_m \leq n, \quad 1 \leq j_1 < j_2 < \dots < j_m \leq n,$$

and exactly c equations $i_k = j_h$ are satisfied. By (4.9), each term in $\sum^{(c)}$ is equal to ξ_c . The number of terms in $\sum^{(c)}$ is

$$\frac{n(n-1) \cdots (n-2m+c+1)}{c!(m-c)!(m-c)!} = \binom{m}{c} \binom{n-m}{m-c} \binom{n}{m} \quad (4.10)$$

Since $\xi_0 = 0$,

$$\text{Var}(U_n) = \binom{n}{m}^{-1} \sum_{c=1}^m \binom{m}{c} \binom{n-m}{m-c} \xi_c \quad (4.11)$$

Hoeffding (1948) obtained the inequality: $0 \leq \xi_c \leq \frac{c}{d} \xi_d \quad 1 \leq c < d \leq m$, which leads to

$$\frac{m^2}{n} \xi_1 \leq \text{Var}(U_n) \leq \frac{m}{n} \xi_m$$

Now, from (4.11) and (4.10)

$$\begin{aligned} \text{Var}(U_n) &= \frac{m^2}{n} \binom{n-m}{n-1} \cdots \binom{n-2m+2}{n-m+1} \xi_1 + \cdots \\ &+ \frac{m!}{n(n-1) \cdots (n-m+1)} \xi_m \end{aligned}$$

Hence $n\text{Var}(U_n)$ is a decreasing function of n which tends to its lower bound $m^2\xi_1$ as n increases, i.e.,

$$\text{Var}(U_n) = \frac{m^2}{n} \xi_1 + O(n^{-2}) \quad (4.12)$$

Therefore, if $E(\phi^2) < \infty$ and $\xi_1 > 0$,

$$n^{1/2}(U_n - \theta) \xrightarrow{d} N(0, m^2 \xi_1), \quad (\text{Hoeffding, 1948}) \quad (4.13)$$

We may rewrite (4.6) as

$$U_n = n^{-[m]} \sum_{1 \leq i_1 \neq \dots \neq i_m \leq n} \int_{R^{pm}} \dots \int \phi(x_1, \dots, x_m) \prod_{j=1}^m d(c(x_j - X_{i_j})),$$

where $n^{-[m]} = (n^{[m]})^{-1} = \{n \dots (n - m + 1)\}^{-1}$.

Writing $d(c(x_j - X_{i_j})) = dF(x_j) + d[c(x_j - X_{i_j}) - F(x_j)]$, $1 \leq j \leq m$, we obtain

$$U_n = \theta(F) + \sum_{h=1}^m \binom{m}{h} U_{n,h} \quad n \geq m \quad (4.14)$$

where

$$U_{n,h} = n^{-[h]} \sum_{1 \leq i_1 \neq \dots \neq i_h \leq n} \int_{R^{ph}} \dots \int \Psi_h(x_1, \dots, x_h) \prod_{j=1}^h d[c(x_j - X_{i_j}) - F(x_j)]$$

for $1 \leq h \leq m$. Further, if we write

$$\begin{aligned} \Psi_h^\circ(x_1, \dots, x_h) &= \Psi_h(x_1, \dots, x_h) - \sum_{j=1}^h \Psi_{h-1}(x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_h) \\ &\quad + \dots + (-1)^h \theta(F), \quad \forall (x_1, \dots, x_h) \in R^{ph}, \end{aligned} \quad (4.15)$$

for $1 \leq h \leq m$, we obtain

$$U_{n,h} = \binom{n}{h}^{-1} \sum_{1 \leq i_1 < \dots < i_h \leq n} \Psi_h^\circ(X_{i_1}, \dots, X_{i_h}), \quad 1 \leq h \leq m \quad (4.16)$$

and the $U_{n,h}$ are themselves U-statistics. From direct computation, $E(U_{n,h}) = 0$, $\forall 1 \leq h \leq m$ and

$$\text{Var}(U_{n,h}) = E(U_{n,h}^2) = O(n^{-h}), \quad h = 1, 2, \dots, m; \quad (4.17)$$

and we can write

$$U_n = \theta(F) + \frac{m}{n} \sum_{i=1}^n [\Psi_1(X_i) - \theta(F)] + O_p(n^{-1}) \quad (4.18)$$

Let $\{\mathbf{X}_i^{(j)}; i \geq 1\}$, $j = 1, \dots, c (\geq 2)$ be independent sequences of independent random vectors, where $\mathbf{X}_i^{(j)}$ has a distribution function $F^{(j)}(\mathbf{x})$, $\mathbf{x} \in R^p$, for $j = 1, \dots, c$. Let $\mathbf{F} = (F^{(1)}, \dots, F^{(c)})$ and $\phi(\mathbf{X}_i^{(j)}, 1 \leq i \leq m_j, 1 \leq j \leq c)$ be a Borel-measurable kernel of degree $\mathbf{m} = (m_1, \dots, m_c)$, where without loss of generality we assume that ϕ is symmetric in the $m_j (\geq 1)$ arguments of the j th set, for $j = 1, \dots, c$. Let $m_0 = m_1 + \dots + m_c$ and

$$\theta(\mathbf{F}) = \int_{R^{m_0}} \dots \int \phi(\mathbf{x}_i^{(j)}, 1 \leq i \leq m_j, 1 \leq j \leq c) \prod_{j=1}^c \prod_{i=1}^{m_j} dF^{(j)}(\mathbf{x}_i^{(j)}) \quad (4.19)$$

Definition 2

For a set of samples of sizes $\mathbf{n} = (n_1, n_2, \dots, n_c)$ with $n_j \geq m_j$, $1 \leq j \leq c$, the *generalized U-statistic* for $\theta(\mathbf{F})$ is

$$U(\mathbf{n}) = \prod_{j=1}^c \binom{n_j}{m_j}^{-1} \sum_{(\mathbf{n})}^* \phi(\mathbf{X}_\alpha^{(j)}, \alpha = i_{j1}, \dots, i_{jm_j}, 1 \leq j \leq c), \quad (4.20)$$

where the summation $\sum_{(\mathbf{n})}^*$ extends over all $1 \leq i_{j1} < \dots < i_{jm_j} \leq n_j$, $1 \leq j \leq c$. $U(\mathbf{n})$ is an unbiased estimator of $\theta(\mathbf{F})$. ■

Now, for every $d_j : 0 \leq d_j \leq m_j$, $1 \leq j \leq c$, let $\mathbf{d} = (d_1, \dots, d_c)$ and

$$\Psi_{d_1 \dots d_c}(\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_{d_j}^{(j)}, 1 \leq j \leq c) \equiv \mathbb{E}(\phi(\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_{d_j}^{(j)}, \mathbf{X}_{d_j+1}^{(j)}, \dots, \mathbf{X}_{m_j}^{(j)}, 1 \leq j \leq c)) \quad (4.21)$$

so that $\Psi_0 = \theta(\mathbf{F})$ and $\Psi_{\mathbf{m}} = \phi$. Then

$$\xi_{\mathbf{d}}(\mathbf{F}) = \mathbb{E} \left(\Psi_{\mathbf{d}}^2(\mathbf{X}_1^{(j)}, \dots, \mathbf{X}_{d_j}^{(j)}, 1 \leq j \leq c) \right) - \theta^2(\mathbf{F}), \quad \mathbf{0} \leq \mathbf{d} \leq \mathbf{m} \quad (4.22)$$

so that $\xi_0(\mathbf{F}) = 0$. Then, for every $\mathbf{n} \geq \mathbf{m}$ (Sen, 1981),

$$\text{Var} [U(\mathbf{n})] = \sum_{j=1}^c n_j^{-1} \sigma_j^2 [1 + O(n_0^{-1})] \quad (4.23)$$

where $n_0 = \min(n_1, \dots, n_c)$ and

$$\sigma_j^2 = m_j^2 \xi_{\delta_{j1}, \dots, \delta_{jc}}(\mathbf{F}) \quad j = 1, \dots, c \quad (4.24)$$

with $\delta_{\alpha\beta} = 1$ or 0 according to whether $\alpha = \beta$ or not.

The decomposition for $U(\mathbf{n})$ can be developed similarly to the one-sample U-statistic. For a two-sample U-statistic of degree (m_1, m_2) , we have

$$\begin{aligned} U(n_1, n_2) &= \theta(\mathbf{F}) + \frac{m_1}{n_1} \sum_{i=1}^{n_1} [\Psi_{10}(X_i) - \theta(\mathbf{F})] + \frac{m_2}{n_2} \sum_{i=1}^{n_2} [\Psi_{01}(Y_i) - \theta(\mathbf{F})] \\ &+ O_p(n_0^{-1}) \end{aligned} \quad (4.25)$$

where $n_0 = \min(n_1, n_2)$.

The above expression can be generalized for multiple-sample U-statistics. For instance, the decomposition for a three-sample and four-sample U-statistics are as follows

$$\begin{aligned} U(n_1, n_2, n_3) &= \theta(\mathbf{F}) + \frac{m_1}{n_1} \sum_{i=1}^{n_1} [\Psi_{100}(X_i) - \theta(\mathbf{F})] + \frac{m_2}{n_2} \sum_{i=1}^{n_2} [\Psi_{010}(Y_i) - \theta(\mathbf{F})] \\ &+ \frac{m_3}{n_3} \sum_{i=1}^{n_3} [\Psi_{001}(Z_i) - \theta(\mathbf{F})] + O_p(n_0^{-1}) \end{aligned} \quad (4.26)$$

where $n_0 = \min(n_1, n_2, n_3)$ and

$$\begin{aligned}
U(n_1, n_2, n_3, n_4) &= \theta(\mathbf{F}) + \frac{m_1}{n_1} \sum_{i=1}^{n_1} [\Psi_{1000}(X_i) - \theta(\mathbf{F})] + \frac{m_2}{n_2} \sum_{i=1}^{n_2} [\Psi_{0100}(Y_i) - \theta(\mathbf{F})] \\
&+ \frac{m_3}{n_3} \sum_{i=1}^{n_3} [\Psi_{0010}(Z_i) - \theta(\mathbf{F})] + \frac{m_4}{n_4} \sum_{i=1}^{n_4} [\Psi_{0001}(W_i) - \theta(\mathbf{F})] \\
&+ O_p(n_0^{-1})
\end{aligned} \tag{4.27}$$

where $n_0 = \min(n_1, n_2, n_3, n_4)$.

4. Combining the U-statistics We can write

$$\begin{aligned}
WSS &= \frac{(N-2)}{3} \sum_{g=1}^G [\mathbf{U}_{1,1}^{(3)} + \mathbf{U}_{1,2}^{(3)} + \mathbf{U}_{1,3}^{(3)}] \\
&+ \frac{(N-2)(N-3)}{12} \sum_{g=1}^G [\mathbf{U}_{2,1}^{(4)} + \mathbf{U}_{2,2}^{(4)} + \mathbf{U}_{2,3}^{(4)}]
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{U}_{1,1}^{(3)} &= \binom{N}{3}^{-1} \sum_{i < j < j'} (D_{ij}^g - D_{ij'}^g)^2, & \mathbf{U}_{1,2}^{(3)} &= \binom{N}{3}^{-1} \sum_{i < i' < j} (D_{ij}^g - D_{i'j}^g)^2 \quad \text{and} \\
\mathbf{U}_{1,3}^{(3)} &= \binom{N}{3}^{-1} \sum_{i < j < j'} (D_{ij}^g - D_{jj'}^g)^2
\end{aligned}$$

are one-sample U-statistics of degree 3 and

$$\begin{aligned}
\mathbf{U}_{2,1}^{(4)} &= \binom{N}{4}^{-1} \sum_{i < j < i' < j'} (D_{ij}^g - D_{i'j'}^g)^2, & \mathbf{U}_{2,2}^{(4)} &= \binom{N}{4}^{-1} \sum_{i < i' < j < j'} (D_{ij}^g - D_{i'j'}^g)^2 \quad \text{and} \\
\mathbf{U}_{2,3}^{(4)} &= \binom{N}{4}^{-1} \sum_{i < i' < j' < j} (D_{ij}^g - D_{i'j'}^g)^2
\end{aligned}$$

are one-sample U-statistics of degree 4. The expected value of WSS is

$$E(WSS) = \sum_{g=1}^G (N-2) \left\{ \mu_{g1} + \frac{(N-3)}{4} \mu_{g2} \right\}.$$

Under H_0 , there is homogeneity among groups, i.e., for any g , $\theta_k^g = \theta_k$ and $\theta_{k_1 k_2}^g = \theta_{k_1 k_2}$, thus

$$E_0(WSS) = G(N-2) \left\{ \mu_1 + \frac{(N-3)}{4} \mu_2 \right\}$$

where

$$\mu_1 = \frac{2}{K^2} \left[\sum_{k=1}^K \theta_k + \sum_{k_1 \neq k_2} \theta_{k_1 k_2} - \sum_{k=1}^K \theta_k(i, j; i, j') - \sum_{k_1 \neq k_2} \theta_{k_1 k_2}(i, j; i, j') \right] \tag{4.28}$$

and

$$\mu_2 = \frac{2}{K^2} \left\{ \sum_{k=1}^K \theta_k (1 - \theta_k) + \sum_{k_1 \neq k_2} (\theta_{k_1 k_2} - \theta_{k_1} \theta_{k_2}) \right\} \quad (4.29)$$

Note that

$$\theta_k = P(X_{ik} \neq X_{jk}) = \sum_{c=0}^{C-1} p_k(c) [1 - p_k(c)] \quad (4.30)$$

$$\begin{aligned} \theta_{k_1 k_2} &= P(X_{ik_1} \neq X_{jk_1}; X_{ik_2} \neq X_{jk_2}) \\ &= \sum_{c_1, c_2=0}^{C-1} p_{k_1 k_2}(c_1, c_2) \left[\sum_{\substack{c_3=0 \\ c_3 \neq c_1}}^{C-1} \sum_{\substack{c_4=0 \\ c_4 \neq c_2}}^{C-1} p_{k_1 k_2}(c_3, c_4) \right] \end{aligned} \quad (4.31)$$

Decomposing WSS , under H_0 ,

$$\begin{aligned} WSS &= G(N-2) \left(\mu_1 + \frac{(N-3)}{4} \mu_2 \right) \\ &+ (N-2) \frac{3}{N} G \sum_{i=1}^N [\Psi_{(1)1}(\mathbf{X}_i) - \mu_1] + O_p(1) \\ &+ \frac{(N-2)(N-3)}{N} G \sum_{i=1}^N [\Psi_{(2)1}(\mathbf{X}_i) - \mu_2] + O_p(N) \end{aligned} \quad (4.32)$$

and the associated mean square expression is

$$\begin{aligned} WMS &\equiv \frac{WSS}{G \binom{N}{2}} = \frac{2WSS}{GN(N-1)} \\ &= \frac{(N-2)(N-3)}{N(N-1)2} \left\{ \mu_2 + \frac{4}{N} \sum_{i=1}^N [\Psi_{(2)1}(\mathbf{X}_i) - \mu_2] \right\} + O_p(N^{-1}) \end{aligned} \quad (4.33)$$

with

$$E_0(WMS) = \frac{\mu_2}{2} + O(N^{-1})$$

and

$$\text{Var}_0(WMS) = \frac{4(N-2)^2(N-3)^2 \xi_1^{(2)}}{GN^2(N-1)^2} + O(N^{-2})$$

For $AWSS$,

$$\begin{aligned} AWSS &= \sum_{1 \leq g_1 < g_2 \leq G} \left[\frac{(N-1)^2}{4} (\mathbf{U}_{4,1}^{(2,2)} + \mathbf{U}_{4,2}^{(2,2)}) \right. \\ &\quad \left. + \frac{(N-1)}{2} (\mathbf{U}_{5,2}^{(2,1)} + \mathbf{U}_{5,1}^{(1,2)}) \right] \end{aligned}$$

where

$$U_{4.1}^{(2,2)} = \left[\binom{N}{2} \binom{N}{2} \right]^{-1} \sum_{\substack{i \neq i' \\ i \neq j'}} \sum_{\substack{j \neq j' \\ j \neq i'}} (D_{ij}^{(g_1, g_2)} - D_{i'j'}^{(g_1, g_2)})^2 \quad \text{and}$$

$$U_{4.2}^{(2,2)} = \left[\binom{N}{2} \binom{N}{2} \right]^{-1} \sum_{\substack{i \neq j \\ i \neq i'}} \sum_{\substack{j \neq j' \\ j \neq i'}} (D_{ij}^{(g_1, g_2)} - D_{i'j'}^{(g_1, g_2)})^2$$

are two-sample U-statistics of degree (2,2) and

$$U_{5.1}^{(1,2)} = \left[\binom{N}{1} \binom{N}{2} \right]^{-1} \sum_{\substack{i=1 \\ i \neq j'}}^N \sum_{\substack{1 \leq j, j' \leq N \\ j \neq j'}} (D_{ij}^{(g_1, g_2)} - D_{ij'}^{(g_1, g_2)})^2 \quad \text{and}$$

$$U_{5.2}^{(2,1)} = \left[\binom{N}{2} \binom{N}{1} \right]^{-1} \sum_{j=1}^N \sum_{\substack{i \neq i' \\ j \neq i'}} (D_{ij}^{(g_1, g_2)} - D_{i'j}^{(g_1, g_2)})^2$$

are two sample U-statistics of degree (1,2) and (2,1), respectively.

$$E(AWSS) = (N-1) \sum_{1 \leq g_1 < g_2 \leq G} \left(\frac{(N-1)}{2} \mu_{(g_1, g_2)4} + \mu_{(g_1, g_2)5} \right)$$

and under H_0

$$E_0(AWSS) = \frac{G(G-1)(N-1)}{2} \left(\frac{(N-1)}{2} \mu_4 + \mu_5 \right)$$

where $\mu_4 = \mu_2$ is given by (4.29) and $\mu_5 = \mu_1$ is given by (4.28).

$AWSS$ can be decomposed as

$$\begin{aligned} AWSS &= \sum_{1 \leq g_1 < g_2 \leq G} \left[\frac{(N-1)^2}{2} \left(\mu_{(g_1, g_2)4} + \frac{2}{N} \sum_{i=1}^N (\Psi_{(4)10}(\mathbf{X}_i^{g_1}) - \mu_{(g_1, g_2)4}) \right. \right. \\ &\quad \left. \left. + \frac{2}{N} \sum_{j=1}^N (\Psi_{(4)01}(\mathbf{X}_j^{g_2}) - \mu_{(g_1, g_2)4}) + O_p(N^{-1}) \right) \right. \\ &\quad \left. + \frac{(N-1)}{2} \left(2\mu_{(g_1, g_2)5} + \frac{1}{N} \sum_{i=1}^N (\Psi_{(5)10}(\mathbf{X}_i^{g_1}) - \mu_{(g_1, g_2)5}) \right. \right. \\ &\quad \left. \left. + \frac{2}{N} \sum_{j=1}^N (\Psi_{(5)01}(\mathbf{X}_j^{g_2}) - \mu_{(g_1, g_2)5}) + \frac{2}{N} \sum_{i=1}^N (\Psi_{(5)10}(\mathbf{X}_i^{g_1}) - \mu_{(g_1, g_2)5}) \right. \right. \\ &\quad \left. \left. + \frac{1}{N} \sum_{j=1}^N (\Psi_{(5)01}(\mathbf{X}_j^{g_2}) - \mu_{(g_1, g_2)5}) + O_p(N^{-1}) \right) \right] \quad (4.34) \end{aligned}$$

The associated mean-square expression is

$$AWMS = \frac{AWSS}{\binom{G}{2} N^2} = \frac{2AWSS}{N^2 G(G-1)}$$

$$\begin{aligned}
&= \frac{(N-1)^2}{N^2 G(G-1)} \sum_{1 \leq g_1 < g_2 \leq G} \left[\mu_{(g_1, g_2)4} + \frac{2}{N} \sum_{i=1}^N (\Psi_{(4)10}(\mathbf{X}_i^{g_1}) - \mu_{(g_1, g_2)4}) \right. \\
&\quad \left. + \frac{2}{N} \sum_{j=1}^N (\Psi_{(4)01}(\mathbf{X}_j^{g_2}) - \mu_{(g_1, g_2)4}) \right] + O_p(N^{-1}) \tag{4.35}
\end{aligned}$$

$$\mathbb{E}_0(AWMS) = \frac{\mu_4}{2} + O(N^{-1})$$

Note that $\mathbb{E}_0(AWMS) = \mathbb{E}_0(WMS)$ since under H_0 , $\mu_4 = \mu_2$,

$$\text{Var}_0(AWMS) = \frac{(N-1)^4}{2N^4 G(G-1)} \left(\frac{4}{N} \xi_{10}^{(4)} + \frac{4}{N} \xi_{01}^{(4)} \right) + O(N^{-2})$$

Now

$$BSS = \frac{N(N-1)}{2} \sum_{g=1}^G (\bar{D}^g - \bar{D}.)^2 = \frac{N(N-1)}{2} \mathbf{D}'_1 \mathbf{D}_1$$

where \mathbf{D}_1 is the $G \times 1$ vector

$$\mathbf{D}_1 = (\bar{D}^1 - \bar{D}., \dots, \bar{D}^G - \bar{D}.)'$$

Note that

$$\mathbb{E}(\bar{D}^g) = \frac{1}{K \binom{N}{2}} \sum_{1 \leq i < j \leq N} \mathbb{E}(D_{ij}^g) = \frac{1}{K} \sum_{k=1}^K \theta_k^g = \bar{\theta}^g$$

$$\mathbb{E}(\bar{D}^{(g_1, g_2)}) = \frac{1}{K} \sum_{k=1}^K \theta_k^{(g_1, g_2)} = \bar{\theta}^{(g_1, g_2)}$$

and

$$\mathbb{E}(\bar{D}.) = \frac{(N-1)}{G(NG-1)} \sum_{g=1}^G \bar{\theta}^g + \frac{2N}{G(NG-1)} \sum_{1 \leq g_1 < g_2 \leq G} \bar{\theta}^{(g_1, g_2)}$$

Therefore,

$$\nu_1 \equiv \mathbb{E}(\bar{D}^g - \bar{D}.) = \bar{\theta}^g - \frac{(N-1)}{G(NG-1)} \sum_{g=1}^G \bar{\theta}^g + \frac{2N}{G(NG-1)} \sum_{1 \leq g_1 < g_2 \leq G} \bar{\theta}^{(g_1, g_2)}$$

Since \bar{D}^g is a U-statistic of degree 2,

$$\text{Var}(\bar{D}^g) = \frac{4}{N} \xi_1^{(12)} + O(N^{-2})$$

where $\xi_1^{(12)} \equiv \mathbb{E}[\psi_{(12)1}^2(\mathbf{X}_i^g)]$, and since $\bar{D}^{(g_1, g_2)}$ is a two-sample U-statistic of degree (1, 1),

$$\text{Var}(\bar{D}^{(g_1, g_2)}) = \frac{1}{N} \xi_{10}^{(13)} + \frac{1}{N} \xi_{01}^{(13)} + O(N^{-2}) \tag{4.36}$$

where $\xi_{10}^{(13)} \equiv \mathbb{E}[\psi_{(13)10}^2(\mathbf{X}_i^{g_1})]$ and $\xi_{01}^{(13)} \equiv \mathbb{E}[\psi_{(13)01}^2(\mathbf{X}_j^{g_2})]$. Under H_0 ,

$$\begin{aligned}\psi_{(12)1}^2(\mathbf{x}_i) &= \frac{1}{K^2} \sum_{k=1}^K \mathbb{P}^2(X_{jk} \neq x_{ik}) + (\bar{\theta})^2 - \frac{2}{K} \bar{\theta} \cdot \sum_{k=1}^K \mathbb{P}(X_{jk} \neq x_{ik}) \\ &+ \frac{1}{K^2} \sum_{k_1 \neq k_2} \mathbb{P}(X_{jk_1} \neq x_{ik_1}; X_{jk_2} \neq x_{ik_2})\end{aligned}$$

We are assuming that under H_0 there is homogeneity across or within groups, i.e., $\theta_k^1 = \theta_k^2 = \dots = \theta_k^G = \theta_k$ and $\theta_k^{(g_1, g_2)} = \theta_k^g = \theta_k$. Therefore, under H_0 ,

$$\sqrt{N} (\bar{D}^g - \bar{\theta}) \xrightarrow{d} \mathbb{N}(0, 4\xi_1^{(12)}) \quad (4.37)$$

and

$$\gamma_{13}^{-1} (\bar{D}^{(g_1, g_2)} - \bar{\theta}) \xrightarrow{d} \mathbb{N}(0, 1) \quad (4.38)$$

where $\gamma_{13}^2 = \frac{1}{N} \xi_{10}^{(13)} + \frac{1}{N} \xi_{01}^{(13)} = \frac{2}{N} \xi_1^{(12)}$ by (4.36).

If $\bar{D}.$ is a linear combination of normal variables, then $\bar{D}.$ also follows a normal distribution.

$$\bar{D} = \frac{(N-1)}{G(NG-1)} \sum_{g=1}^G \bar{D}^g + \frac{2N}{G(NG-1)} \sum_{1 \leq g_1 < g_2 \leq G} \bar{D}^{(g_1, g_2)}$$

Under H_0 ,

$$\eta_1 \equiv \mathbb{E}_0(\bar{D}.) = \frac{(N-1)\bar{\theta} + N(G-1)\bar{\theta}}{(NG-1)} = \bar{\theta}.$$

$$\begin{aligned}\sigma_1^2 &\equiv \text{Var}_0(\bar{D}.) \\ &+ \frac{4N^2}{G^2(NG-1)^2} \\ &= \frac{(N-1)^2}{G(NG-1)^2} \frac{4}{N} \xi_1^{(12)} + \frac{2N^2(G-1)}{G(NG-1)^2} \left[\left(\frac{1}{N} (\xi_{10}^{(13)} + \xi_{01}^{(13)}) \right) \right. \\ &\left. + 2(G-2) \frac{1}{N} \xi_{10}^{(13,1;13,2)} \right] + \frac{2N(N-1)}{G^2(NG-1)^2} G(G-1) \frac{2}{N} \xi_1^{(12,13)}\end{aligned}$$

where $\xi_{10}^{(13,1;13,2)} = \mathbb{E}\{\psi_{(13,1)10}(\mathbf{X}_i^{g_1}) \psi_{(13,2)10}(\mathbf{X}_i^{g_1})\}$ and

$\psi_{(13,2)10}(\mathbf{x}_i^{g_1}) = \mathbb{E}[\phi_{13,2}(\mathbf{x}_i^{g_1}, \mathbf{X}_j^{g_3}) - \bar{\theta}^{(g_1, g_3)}]$. Under H_0 ,

$\psi_{(12)1}(\mathbf{X}_i) = \psi_{(13)10}(\mathbf{X}_i) = \psi_{(13)01}(\mathbf{X}_j) = \psi_{(13,1)10}(\mathbf{X}_i) = \psi_{(13,2)10}(\mathbf{X}_i)$. Therefore, $\xi_1^{(12)} = \xi_{10}^{(13)} = \xi_{01}^{(31)} = \xi_{10}^{(13,1;13,2)} = \xi_1^{(12,13)}$ and

$$\sigma_1^2 = [(N-1)^2 + N(G-1)(NG-1)] \frac{4\xi_1^{(12)}}{NG(NG-1)^2} \quad (4.39)$$

Hence, under H_0 ,

$$\sigma_1^{-1} (\bar{D} - \bar{\theta}) \xrightarrow{d} \mathbb{N}(0, 1)$$

Now

$$\nu_1 = \mathbf{E}_0(\bar{D}^g - \bar{D}) = \bar{\theta} - \bar{\theta} = 0 \quad (4.40)$$

and

$$\begin{aligned} \tau_1^2 &\equiv \text{Var}_0(\bar{D}^g - \bar{D}) \\ &= \left[1 - 2 \frac{(N-1)}{G(NG-1)} \right] \frac{4}{N} \xi_1^{(12)} + \sigma_1^2 - \frac{4N(G-1)}{G(NG-1)} \frac{2}{N} \xi_1^{(12,13)} \end{aligned} \quad (4.41)$$

where $\xi_1^{(12,13)} \equiv \mathbf{E}\{\psi_{(12)1}(\mathbf{X}_i^{g_1})\psi_{(13)10}(\mathbf{X}_i^{g_1})\} = \xi_1^{(12)}$, since $\psi_{(12)1}(\mathbf{X}_i) = \psi_{(13)10}(\mathbf{X}_i)$ under H_0 .

Then,

$$\tau_1^2 = \{(N-1)^2 + (NG-1)[N(G-1) + (NG-1)(G-2)]\} \frac{4\xi_1^{(12)}}{NG(NG-1)^2} \quad (4.42)$$

So,

$$\tau_1^{-1}(\bar{D}^g - \bar{D}) \xrightarrow{d} \mathbf{N}(0, 1)$$

Since BSS is a quadratic form of normal random variables,

$$BSS = \frac{N(N-1)}{2} \mathbf{D}'_1 \mathbf{D}_1 \sim \frac{N(N-1)}{2} \sum_{g=1}^G \lambda_g (\chi_1^2)_g$$

which is a linear combination of χ_1^2 random variables, where λ_g 's are the characteristic roots of $\text{Var}(\mathbf{D}_1) = \boldsymbol{\Sigma}_1$. Note that the diagonal elements of $\boldsymbol{\Sigma}_1$ are τ_1^2 and the off-diagonal elements, under H_0 , are

$$\begin{aligned} \text{Cov}_0(\bar{D}^{g_1} - \bar{D}, \bar{D}^{g_2} - \bar{D}) \\ = \left[\frac{(N-1)^2 - (NG-1)(NG+N-2)}{(NG-1)} \right] \frac{4\xi_1^{(12)}}{NG(NG-1)} < 0 \end{aligned}$$

since $(NG-1)(NG+N-2) > (N-1)^2$.

Now,

$$\mathbf{E}_0(BSS) = \frac{N(N-1)}{2} \text{trace}(\boldsymbol{\Sigma}_1) = \frac{N(N-1)}{2} G \tau_1^2$$

and

$$\text{Var}_0(BSS) = \frac{N^2(N-1)^2}{4} \text{trace}(\boldsymbol{\Sigma}_1)^2$$

Let

$$BMS = \frac{BSS}{G \binom{N}{2}} = \frac{1}{G} \mathbf{D}'_1 \mathbf{D}_1$$

Then

$$\mathbf{E}_0(BMS) = \frac{1}{G} \mathbf{E}_0(BSS) = \tau_1^2$$

and

$$\text{Var}_0(BMS) = \frac{1}{G^2} \text{Var}_0(BSS) = \frac{1}{G^2} \text{trace}(\boldsymbol{\Sigma}_1)^2$$

For *ABSS* we have,

$$ABSS = \sum_{1 \leq g_1 < g_2 \leq G} \sum_{i=1}^N \sum_{j=1}^N (\bar{D}^{(g_1, g_2)} - \bar{D}.)^2 = N^2 \mathbf{D}_2 \mathbf{D}_2$$

where $\mathbf{D}_2 = (\bar{D}^{(1,2)} - \bar{D}., \bar{D}^{(1,3)} - \bar{D}., \dots, \bar{D}^{(G-1,G)} - \bar{D}.)'$ is a $\frac{G(G-1)}{2} \times 1$ vector.

Let

$$\nu_2 \equiv \mathbf{E}(\bar{D}^{(g_1, g_2)} - \bar{D}.) = \bar{\theta}^{(g_1, g_2)} - \frac{(N-1)}{G(NG-1)} \sum_{g=1}^G \bar{\theta}^g - \frac{2N}{G(NG-1)} \sum_{1 \leq g_1 < g_2 \leq G} \bar{\theta}^{(g_1, g_2)}$$

Under H_0 ,

$$\nu_2 = \mathbf{E}_0(\bar{D}^{(g_1, g_2)} - \bar{D}.) = \bar{\theta} - \bar{\theta} = 0 \quad (4.43)$$

and

$$\begin{aligned} \tau_2^2 &\equiv \text{Var}(\bar{D}^{(g_1, g_2)} - \bar{D}.) \\ &= \text{Var}(\bar{D}^{(g_1, g_2)}) + \text{Var}(\bar{D}.) - 2\text{Cov}(\bar{D}^{(g_1, g_2)}, \bar{D}.) \\ &= \frac{1}{N} (\xi_{10}^{(13)} + \xi_{01}^{(13)}) + \sigma_1^2 - \frac{4(N-1)}{G(NG-1)} \frac{2}{N} \xi_1^{(12,13)} \\ &\quad - \frac{4N}{G(NG-1)} \left[\frac{1}{N} (\xi_{10}^{(13)} + \xi_{01}^{(13)}) + 2(G-2) \frac{1}{N} \xi_{10}^{(13,1;13,2)} \right] \end{aligned} \quad (4.44)$$

Note that under H_0 there is homogeneity among groups,

$$\Psi_{(13)10}(\mathbf{x}_i) = \Psi_{(13)01}(\mathbf{x}_j) = \Psi_{(13,1)10}(\mathbf{x}_i) = \Psi_{(13,2)10}(\mathbf{x}_i) = \frac{1}{K} \sum_{k=1}^K \mathbf{P}(X_{ik} \neq x_{jk})$$

since the sequences are i.i.d.

Therefore, $\Psi_{(13,1)10}(\mathbf{x}_i) \Psi_{(13,2)10}(\mathbf{x}_i) = \Psi_{(13)10}^2(\mathbf{x}_i)$ and

$$\xi_{10}^{(13,1;13,2)} = \xi_{10}^{(13)} = \xi_{01}^{(13)} = \xi_1^{(12)} = \xi_1^{(12,13)}$$

So, under H_0 ,

$$\tau_2^2 = \{2(N-1)^2 + (NG-1)[2N(G-1) + (NG-1)(G-4)]\} \frac{2\xi_1^{(12)}}{NG(NG-1)^2} \quad (4.45)$$

As in *BSS*,

$$ABSS \sim N^2 \sum_{i=1}^{G(G-1)/2} \lambda_i (\chi_1^2)_i$$

where λ_i 's are the characteristic roots of $\boldsymbol{\Sigma}_2 = \text{Var}(\mathbf{D}_2)$. The diagonal elements of $\boldsymbol{\Sigma}_2$ are τ_2^2 and, if all groups are different, the off-diagonal elements are

$$\begin{aligned} \text{Cov}(\bar{D}^{(g_1, g_2)} - \bar{D}, \bar{D}^{(g_3, g_4)} - \bar{D}) \\ = [(N-1)^2 - (NG-1)(NG+N-2)] \frac{4\xi_1^{(12)}}{NG(NG-1)^2} < 0 \end{aligned}$$

and if $g_1 = g_2$ or $g_1 = g_3$ or $g_2 = g_3$,

$$\begin{aligned} \text{Cov}(\bar{D}^{(g_1, g_2)} - \bar{D}, \bar{D}^{(g_1, g_3)} - \bar{D}) \\ = \{4(N-1)^2 + (NG-1)[4N(G-1) + (G-8)(NG-1)]\} \frac{\xi_1^{(12)}}{NG(NG-1)^2} . \end{aligned}$$

Now

$$E_0(ABSS) = N^2 \text{trace}(\boldsymbol{\Sigma}_2) = N^2 \frac{G(G-1)}{2} \tau_2^2$$

$$\text{Var}_0(ABSS) = N^4 \text{trace}(\boldsymbol{\Sigma}_2)^2$$

The corresponding mean-square term is defined as

$$ABMS = \frac{ABSS}{N^2 \binom{G}{2}} = \frac{2}{G(G-1)} \mathbf{D}'_2 \mathbf{D}_2$$

Then

$$E_0(ABMS) = \frac{2}{G(G-1)} \text{trace}(\boldsymbol{\Sigma}_2) = \tau_2^2$$

$$\text{Var}_0(ABMS) = \frac{4}{G^2(G-1)^2} \text{trace}(\boldsymbol{\Sigma}_2)^2$$

5. Test Statistics One alternative is to compare *WMS* with *AWMS*. Let $T_1 = \frac{WMS}{AWMS}$. Under H_0 ,

$$\frac{WMS}{AWMS} = \frac{\frac{(N-2)(N-3)}{2N(N-1)} \left\{ \mu_2 + \frac{4}{N} \sum_{i=1}^N (\Psi_{(2)1}(\mathbf{X}_i) - \mu_2) \right\} + O_p(N^{-1})}{\frac{(N-1)^2}{2N^2} \left\{ \mu_2 + \frac{4}{N} \sum_{i=1}^N (\Psi_{(2)1}(\mathbf{X}_i) - \mu_2) \right\} + O_p(N^{-1})}$$

But, $\frac{WMS}{AWMS} \xrightarrow{p} 1$ as $N \rightarrow \infty$, i.e, asymptotically the distribution of $\frac{WMS}{AWMS}$ is degenerate.

Let $\Sigma_1 = \frac{1}{N}\Sigma_1^*$ and $\Sigma_2 = \frac{1}{N}\Sigma_2^*$. Under H_0 ,

$$BMS = \frac{BSS}{G\binom{N}{2}} \sim \frac{1}{NG} \sum_{g=1}^G \lambda_{1g}^* (\chi_{1g}^2)$$

$$ABMS = \frac{ABSS}{N^2\binom{G}{2}} \sim \frac{2}{NG(G-1)} \sum_{i=1}^{G(G-1)/2} \lambda_{2i}^* (\chi_{1i}^2)$$

where λ_{1g}^* 's and λ_{2i}^* 's are the characteristic roots of Σ_1^* and Σ_2^* , respectively. Also, under H_0 , by theoretical results pertaining to U-statistics

$$\sqrt{N}(WMS - \mu_2/2) \rightarrow N\left(0, \frac{4}{G}\xi_1^{(2)}\right)$$

and

$$\sqrt{N}(AWMS - \mu_2/2) \rightarrow N\left(0, \frac{4}{G(G-1)}\xi_1^{(2)}\right).$$

Thus,

$$BMS = O_p(N^{-1}) \quad \text{and} \quad ABMS = O_p(N^{-1})$$

while

$$WMS = O_p(N^{-1/2}) \quad \text{and} \quad AWMS = O_p(N^{-1/2})$$

Define

$$T_{N,2} \equiv N\left(\frac{BMS}{WMS}\right) \quad \text{and} \quad T_{N,3} \equiv N\left(\frac{ABMS}{AWMS}\right).$$

Since, BMS and $ABMS$ are the dominating terms in $T_{N,2}$ and $T_{N,3}$, respectively, we can write

$$T_{N,2} = \frac{2N(BMS)}{\mu_2} + O_p(N^{-1/2})$$

and

$$T_{N,3} = \frac{2N(ABMS)}{\mu_2} + O_p(N^{-1/2})$$

Therefore,

$$T_{N,2} \sim \frac{2}{G\mu_2} \sum_{g=1}^G \lambda_{1g}^* (\chi_{1g}^2)$$

and

$$T_{N,3} \sim \frac{4}{G(G-1)\mu_2} \sum_{i=1}^{G(G-1)/2} \lambda_{2i}^* (\chi_{1i}^2)$$

Because the elements of $\boldsymbol{\Sigma}_1^*$ and $\boldsymbol{\Sigma}_2^*$ are unknown, the characteristic roots of these matrices are also unknown. Therefore, the above distributions do not have a closed analytic form and we call upon resampling methods, such as the bootstrap, to generate the reference distribution for the test statistic.

6. Power of the Tests

Lemma 1

Let \mathbf{T}_n be a vector of random variables that can be expressed as

$$\mathbf{T}_n = \boldsymbol{\nu} + \frac{1}{\sqrt{n}}\mathbf{U}_n + \mathbf{R}_n$$

where $\mathbf{R}_n = O_p(n^{-1})$.

If $Q(\mathbf{T}) = \mathbf{T}'\mathbf{A}\mathbf{T}$ is a quadratic form on \mathbf{T} . Then,

$$\begin{aligned} Q(\mathbf{T}) &= \mathbf{T}'\mathbf{A}\mathbf{T} = \left\{ \boldsymbol{\nu} + \frac{1}{\sqrt{n}}\mathbf{U}_n + \mathbf{R}_n \right\}' \mathbf{A} \left\{ \boldsymbol{\nu} + \frac{1}{\sqrt{n}}\mathbf{U}_n + \mathbf{R}_n \right\} \\ &= Q(\boldsymbol{\nu}) + \frac{2}{\sqrt{n}}\boldsymbol{\nu}'\mathbf{A}\mathbf{U}_n + \frac{1}{n}Q(\mathbf{U}_n) + 2\boldsymbol{\nu}'\mathbf{A}\mathbf{R}_n + O_p(n^{-3/2}) \end{aligned}$$

If $\boldsymbol{\nu} = \mathbf{0}$ then $Q(\mathbf{T}) = \frac{1}{n}Q(\mathbf{U}_n) + O_p(n^{-3/2})$.

In our case, $\mathbf{T} = \mathbf{D}_1$ and the quadratic form is $Q(\mathbf{D}_1) = \mathbf{D}_1'\mathbf{D}_1$. Note that we can write,

$$\mathbf{D}_1'\mathbf{D}_1 = \sum_{g=1}^G (\bar{D}^g - \bar{D}.)^2 = \sum_{g=1}^G (\bar{D}^g - \bar{D}.)^2 + 2\nu_1 \sum_{g=1}^G (\bar{D}^g - \bar{D}.) + G\nu_1^2$$

Let $\mathbf{V}_N = \mathbf{D}_1 - \boldsymbol{\nu}_1$, where $\boldsymbol{\nu}_1$ is a vector $G \times 1$ with elements ν_1 . Then, $E(\mathbf{V}_N) = \mathbf{0}$ and $\text{Var}(\mathbf{V}_N) = \boldsymbol{\Sigma}_1 = \frac{1}{N}\boldsymbol{\Sigma}_1^* = O(N^{-1})$. Therefore,

$$Q(\mathbf{D}_1) = \mathbf{D}_1'\mathbf{D}_1 = \mathbf{V}_N'\mathbf{V}_N + 2\boldsymbol{\nu}_1'\mathbf{V}_N + \boldsymbol{\nu}_1'\boldsymbol{\nu}_1$$

Since $\sqrt{N}\mathbf{V}_N \sim N(\mathbf{0}, \boldsymbol{\Sigma}_1^*)$,

$$N\mathbf{V}_N'\mathbf{V}_N \sim \sum_{g=1}^G \lambda_g^* (\chi_1^2)_g$$

where λ_g^* are the characteristic roots of $\boldsymbol{\Sigma}_1^*$. Also,

$$2\sqrt{N}\boldsymbol{\nu}_1'\mathbf{V}_N \sim N(\mathbf{0}, 4\boldsymbol{\nu}_1'\boldsymbol{\Sigma}_1^*\boldsymbol{\nu}_1)$$

Now,

$$T_{N,2} = \frac{2N}{G\mu_2}\mathbf{V}_N'\mathbf{V}_N + \frac{4\sqrt{N}\boldsymbol{\nu}_1'}{G\mu_2} \left(\sqrt{N}\mathbf{V}_N \right) + \frac{2N}{\mu_2}\nu_1^2 + O_p(N^{-1/2})$$

$$\left(\frac{T_{N,2} - 2N\nu_1^2/\mu_2}{4\sqrt{N}\nu_1/(G\mu_2)} \right) = \frac{N \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1)^2}{2\sqrt{N}\nu_1} + \sqrt{N} \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1) + O_p(N^{-1})$$

Note that

$$\frac{N \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1)^2}{2\sqrt{N}\nu_1} = O_p(N^{-1/2}), \text{ since } N \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1)^2 = O_p(1)$$

and

$$\sqrt{N} \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1) = O_p(1), \text{ since } \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1) = O_p(N^{-1/2})$$

So, for a fixed $\nu_1 \neq 0$, as $N \rightarrow \infty$,

$$\left(\frac{T_{N,2} - 2N\nu_1^2/\mu_2}{4\sqrt{N}\nu_1/(G\mu_2)} \right) = \sqrt{N} \sum_{g=1}^G (\bar{D}^g - \bar{D} - \nu_1) + O_p(N^{-1/2})$$

Thus,

$$P(T_{N,2} > \nu_1) = P\left(Z > G \frac{(\mu_2 - 2N\nu_1)}{4\sqrt{N}} \right) \rightarrow 1, \text{ as } N \rightarrow \infty,$$

i.e., this test is consistent.

Now, consider a local alternative hypothesis. Let $\nu_1 = \frac{1}{\sqrt{N}}\gamma_1^*$, where γ_1^* is a constant. Then,

$$\begin{aligned} T_{N,2} &= \frac{2N}{G\mu_2} \mathbf{V}'_N \mathbf{V}_N + \frac{4\gamma_1^*}{G\mu_2} \left[\sqrt{N} \sum_{g=1}^G \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N}}\gamma_1^* \right) \right] \\ &\quad + \frac{2}{\mu_2} (\gamma_1^*)^2 + O_p(N^{-1/2}) \end{aligned}$$

$$\begin{aligned} \left(\frac{T_{N,2} - 2(\gamma_1^*)^2/\mu_2}{4\gamma_1^*/(G\mu_2)} \right) &= \frac{N \sum_{g=1}^G \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N}}\gamma_1^* \right)^2}{2\gamma_1^*} \\ &\quad + \sqrt{N} \sum_{g=1}^G \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N}}\gamma_1^* \right) + O_p(N^{-1/2}) \end{aligned}$$

Note that

$$\frac{N \sum_{g=1}^G \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N}}\gamma_1^* \right)^2}{2\gamma_1^*} = O_p(1) \quad \text{and} \quad \sqrt{N} \sum_{g=1}^G \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N}}\gamma_1^* \right) = O_p(1)$$

Therefore, $T_{N,2}$ no longer follows a Normal distribution as $N \rightarrow \infty$. It is a convolution of a linear combination of chi-square random variables and a normal random variable:

$$T_{N,2} = \frac{2N}{G\mu_2} \mathbf{V}'_N \mathbf{V}_N + \frac{4\sqrt{N}}{G\mu_2} (\gamma_1^*)' \mathbf{V}_N + \frac{2(\gamma_1^*)^2}{\mu_2} + O_p(N^{-1/2})$$

$$T_{N,2} \sim \frac{2}{G\mu_2} \sum_{g=1}^G \lambda_{1g}^* (\chi_1^2)_g + N \left(\mathbf{0}, \frac{16}{G^2\mu_2^2} (\gamma_1^*)' \boldsymbol{\Sigma}_1^* \gamma_1^* \right) + \frac{2(\gamma_1^*)^2}{\mu_2}$$

Now, let us find out whether $\mathbf{V}'_N \mathbf{V}_N$ and $(\gamma_1^*)' \mathbf{V}_N$ are independent. $\mathbf{V}'_N \mathbf{V}_N$ and $(\gamma_1^*)' \mathbf{V}_N$ are independent if and only if $(\gamma_1^*)' \boldsymbol{\Sigma}_1 = \mathbf{0}$ (Searle, 1971).

Recall that

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} \tau_1^2 & \tau_{12} & \dots & \tau_{12} \\ \tau_{12} & \tau_1^2 & \dots & \tau_{12} \\ \vdots & \vdots & \ddots & \vdots \\ \tau_{12} & \tau_{12} & \dots & \tau_1^2 \end{pmatrix}$$

where

$$\tau_1^2 = \{(N-1)^2 + (NG-1)[N(G-1) + (NG-1)(G-2)]\} \frac{4\xi_1^{(12)}}{NG(NG-1)^2}$$

and

$$\tau_{12} = \{(N-1)^2 - (NG-1)(NG+N-2)\} \frac{4\xi_1^{(12)}}{NG(NG-1)^2}$$

Then,

$$(\gamma_1^*)' \boldsymbol{\Sigma}_1 = \gamma_1^* [\tau_1^2 + (G-1)\tau_{12} \dots \tau_1^2 + (G-1)\tau_{12}]$$

and

$$\begin{aligned} \tau_1^2 + (G-1)\tau_{12} &= 0 \\ \Leftrightarrow G(N-1)^2 + (NG-1)[N(G-1) + (NG-1)(G-2)] \\ &\quad - (G-1)(NG+N-2) = 0 \\ \Leftrightarrow N &= 1 \end{aligned}$$

So, $\mathbf{V}'_N \mathbf{V}_N$ and $(\gamma_1^*)' \mathbf{V}_N$ are independent if and only if $N = 1$, which is not the case here.

Now, write

$$\frac{2}{G\mu_2} [N\mathbf{V}'_N \mathbf{V}_N + 2\sqrt{N}(\gamma_1^*)' \mathbf{V}_N] = \frac{2}{G\mu_2} [(\sqrt{N}\mathbf{V}_N + \gamma_1^*)'(\sqrt{N}\mathbf{V}_N + \gamma_1^*) - (\gamma_1^*)' \gamma_1^*]$$

and

$$\begin{aligned} T_{N,2} &= \frac{2N}{\mu_2} (BMS) = \frac{2N}{G\mu_2} \mathbf{D}'_1 \mathbf{D}_1 \\ &= \frac{2}{G\mu_2} (\sqrt{N}\mathbf{V}_N + \gamma_1^*)' (\sqrt{N}\mathbf{V}_N + \gamma_1^*) + O_p(N^{-1/2}) \end{aligned}$$

Note that $\sqrt{N}\mathbf{V}_N + \boldsymbol{\gamma}_1^* \sim N(\boldsymbol{\gamma}_1^*, \boldsymbol{\Sigma}_1^*)$ and

$$\mathbf{D}_1 \sim N(\boldsymbol{\nu}_1, \boldsymbol{\Sigma}_1) \quad \text{or} \quad \sqrt{N}\mathbf{D}_1 \sim N(\boldsymbol{\gamma}_1^*, \boldsymbol{\Sigma}_1^*).$$

The distribution of $\sqrt{N}\mathbf{D}_1'\mathbf{D}_1$ can also be derived the following way.

Let \mathbf{P} be a $G \times G$ orthogonal matrix (i.e., $\mathbf{P}'\mathbf{P} = \mathbf{I}$) such that $\mathbf{P}\boldsymbol{\Sigma}_1^*\mathbf{P}' = \boldsymbol{\Lambda}$, where $\boldsymbol{\Lambda}$ is a diagonal matrix, and

$$\mathbf{Y} = \sqrt{N}\mathbf{P}\mathbf{D}_1 \Rightarrow \sqrt{N}\mathbf{D}_1 = \mathbf{P}'\mathbf{Y}$$

Then,

$$\mathbf{Y} \sim N(\mathbf{P}\boldsymbol{\gamma}_1^*, \boldsymbol{\Lambda}) \quad \text{and} \quad N\mathbf{D}_1'\mathbf{D}_1 = \mathbf{Y}'\mathbf{P}\mathbf{P}'\mathbf{Y} = \mathbf{Y}'\mathbf{Y},$$

Hence,

$$N\mathbf{D}_1'\mathbf{D}_1 = \mathbf{Y}'\mathbf{Y} \sim \sum_{i=1}^G \lambda_i (\chi_1^2(\delta_i)) \quad (6.46)$$

where $\delta_i = \frac{(\nu_{1i}^*)^2}{\lambda_i}$, λ_i 's are the diagonal elements of the diagonal matrix $\boldsymbol{\Lambda}$ and ν_{1i}^* is the i th row of the vector $\boldsymbol{\nu}_1^* = \mathbf{P}\boldsymbol{\gamma}_1^*$. By (6.46),

$$T_{N,2} = \frac{2N}{G\mu_2}\mathbf{D}_1'\mathbf{D}_1 \sim \frac{2}{G\mu_2} \sum_{i=1}^G \lambda_i (\chi_1^2(\delta_i))_i$$

Since we have a linear combination of non-central chi-square random variables, when $\nu_1 = \frac{\gamma_1^*}{\sqrt{N}}$,

$$P(T_{N,2} > \nu_1) \rightarrow 1 \quad \text{as} \quad N \rightarrow \infty$$

As the distribution of $T_{N,3}$ is similar to the distribution of $T_{N,2}$, the above results about consistency and power of the test apply to $T_{N,3}$.

7. Data analysis The data set consists of two groups of HIV infected individuals (subtype B and not B) with 46 sequences (individuals) each. The nucleotide sequences are all from epidemiologically independent individuals, which means that among the individuals in our sample, one did not infect the other, i.e., they were not sharing the same syringe, they were not partners. Since the sequences are in the nucleotide level, there are therefore four categories. After aligning the sequences and discarding the positions with no change, we end up with 155 positions.

Since the elements of $\boldsymbol{\Sigma}_1^*$ and $\boldsymbol{\Sigma}_2^*$ are unknown, the characteristic roots of these matrices are also unknown and the distributions of the test statistics do not have a closed analytic form. In view of this, we call upon resampling techniques, such as the bootstrap. Here is a summary of the procedure:

1. Compute the statistics T_{N2} and T_{N3} from the data set.
2. Sample 46 sequences to each group with replacement from the pooled sample, i.e., the combined groups.
3. Recompute the test statistics T_{N2} and T_{N3} from this sample and store it.
4. Repeat steps 2 and 3 R times (R should be at least 1,000).

The p-values for the tests are then $\frac{\#T'_{N2s} \geq T_{N2obs}}{R}$ and $\frac{\#T'_{N3s} \geq T_{N3obs}}{R}$.

The results are

$$T_{N2obs} = 20,07 \quad T_{N3obs} = 4,17$$

For $R = 10,000$, the percentiles of the bootstrap distribution are given in Table 1 and the observed p-value for T_{N2} and T_{N3} are less than $1/10001$. Therefore, we can say that relative to the within-clade variation, there is significant variability between the two clades and similarly, relative to the across-within-clade, there is significant variability across-between the two clades.

Table 1: Percentiles of the Bootstrap Distribution

Statistic	1%	5%	95%	99%
T_{N2}	0.0002	0.0007	2.6799	4.1866
T_{N3}	0.0000	0.0000	0.0132	0.0520

References

- [1] R J Anderson and J R Landis. CATANOVA for multidimensional contingency tables: Nominal-scale response. *Communications in Statistics - Theory and Methods*, A9(11):1191–1206, 1980.
- [2] R J Anderson and J R Landis. CATANOVA for multidimensional contingency tables: Ordinal-scale response. *Communications in Statistics - Theory and Methods*, 11(3):257–270, 1982.
- [3] W Hoeffding. A class of statistics with asymptotically normal distribution. *Ann. of Math. Stat.*, 19:293–325, 1948.
- [4] A J Lee. *U-Statistics - Theory and Practice*. Marcel Dekker, Inc, 1990.
- [5] R J Light and B H Margolin. An analysis of variance for categorical data. *Journal of the American Statistical Association*, 66:534–544, 1971.

- [6] R J Light and B H Margolin. An analysis of variance for categorical data II: Small sample comparisons with chi square and other competitors. *Journal of the American Statistical Association*, 69:755–764, 1974.
- [7] R V Mises. On the asymptotic distribution of the differentiable statistical functions. *Annals of Mathematical Statistics*, 18:309–348, 1947.
- [8] H P Pinheiro. *Modelling Variability in the HIV Genome*. PhD thesis, University of North Carolina, December 1997. Mimeo Series No. 2186T.
- [9] H P Pinheiro, F Seillier-Moiseiwitsch, and P K Sen. Multivariate CATANOVA and applications to DNA sequences in categorical data. Research Report 12/99, Universidade Estadual de Campinas, 1999.
- [10] M L Puri and P K Sen. *Nonparametric Methods in Multivariate Analysis*. Wiley, New York, 1971.
- [11] S R Searle. *Linear Models*. John Wiley & Sons, 1971.
- [12] S R Searle. *Matrix Algebra Useful for Statistics*. John Wiley & Sons, 1982.
- [13] F Seillier-Moiseiwitsch, B H Margolin, and R Swanstrom. Genetic variability of human immunodeficiency virus: Statistical and biological issues. *Annual Review of Genetics*, 28:559–596, 1994.
- [14] P K Sen. On some multisample permutation tests based on a class of u-statistics. *American Statistical Association Journal*, pages 1201–1213, 1967.
- [15] P K Sen. *Invariance Principles and Statistical Inference*. John Wiley & Sons, 1981.
- [16] P K Sen. Paired comparisons for multiple characteristics: An anocova approach. In H N Nagaraja, P K Sen, and D F Morrison, editors, *Statistical Theory and Practice: Papers in Honor of H. A. David*, pages 247–264. Springer-Verlag, New York, 1995.
- [17] P K Sen and J M Singer. *Large Sample Methods in Statistics*. Chapman & Hall, 1993.
- [18] B S Weir. *Genetic Data Analysis*. Sinauer Associates, 1990.