

Analysis of Variance for Hamming Distances Applied to Unbalanced Designs *

BY HILDETE PRISCO PINHEIRO

Department of Statistics

State University of Campinas, Brazil

FRANÇOISE SEILLIER-MOISEIWITSCH

Department of Biostatistics

University of North Carolina at Chapel Hill, U.S.A.

PRANAB KUMAR SEN

Department of Biostatistics

University of North Carolina at Chapel Hill

Abstract

The interest here is the between- and within-group comparison of genomic sequences. All possible pairwise comparisons within and across groups are performed. Thus, unlike in analyses relying on measures of diversity (such as the Gini-Simpson index), sequences are considered on an individual basis. We develop a categorical analysis-of-variance framework for Hamming distances. This metric measures the proportion of positions at which two aligned sequences differ. We assume that the sequences are distantly related, but do not require that positions along the genome be independent. The total sum of squares is decomposed into within-, between- and across-group expressions. The latter term does not appear in the classical set-up. The theory of generalized U-statistics is utilized to find the asymptotic distribution of each sum of squares. Test statistics to assess homogeneity among groups are constructed.

1. Introduction

The focus of this paper lies in the comparison of genomic sequences, be they either DNA or protein sequences. These sequences are grouped in some way and the heterogeneity in these groups is quantified. The objective is to assess whether the variability is constant across groups. For instance, for sequences from the human immunodeficiency virus (HIV), a number of papers have investigated whether heterogeneity in the envelope gene is the same for all clades (see Seillier-Moiseiwitsch et al. (1994) for a review).

*This research was funded in part by Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (202/92-02), Fundação de Amparo à Pesquisa do Estado de São Paulo (1998/12199-4), the National Science Foundation (DMS-9305588), the American Foundation for AIDS Research (70428-15-RF) and the National Institutes of Health (R29-GM49804 and P30-HD37260).

Key words and phrases: Amino Acid, Analysis of Variance, Categorical Data, Genome, Hamming Distance, Nucleotide, U-statistics, Bootstrap

Weir (1990) described an analysis of variance for the genetic variation in populations, as measured by the observed *heterozygosity*. The variance of the average heterozygosity is broken down to show the contribution of populations, loci and individuals by setting out the calculations as in an analysis of variance. Our situation is somewhat different because we would like to consider genomic regions rather than a small number of loci. This is why we selected the Hamming distance as our metric. The Hamming distance is the proportion of positions at which two aligned sequences differ. The sequences are regarded as independent but no assumption is imposed on the correlation between positions.

We develop an analysis-of-variance framework based on Hamming distances. 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. We estimate the variability between, within and across groups (Section 2). In the within-group sum of squares, we are estimating the variability among sequences within a group around the average distance within this group. In the across-group sum of squares, we are estimating the variability of sequences across two groups with respect to the average distance between those groups. In the between-group sum of squares, we estimate the variability in the group average distances around the overall average. U-statistics are utilized to represent the average distance between and within groups as well as the overall distance (Sections 3, 4 and 5). The total sum of squares is decomposed into within-, between- and across-group sums of squares. The latter term does not appear in the classical set-up. The theory for generalized U-statistics (Puri & Sen, 1971; Lee, 1990; Sen & Singer, 1993) is used to find the asymptotic distributions of these sums of squares. In Section 6, test statistics are developed to assess homogeneity among groups. The power of the tests is discussed in Section 7. A data analysis is described briefly in Section 8. The paper closes with a discussion in Section 9.

2. The Total Sum of Squares and its Decomposition

Let X_{ik}^g be the label (i.e., an amino acid for protein sequences or a nucleotide $x_{ik}^g \in \{\mathbf{A}, \mathbf{C}, \mathbf{T}, \mathbf{G}\}$ for DNA sequences) present at position k ($k = 1, \dots, K$) of sequence i ($i = 1, \dots, N_g$) in group g ($g = 1, \dots, G$). Then $\mathbf{X}_i^g = (X_{i1}^g, X_{i2}^g, \dots, X_{iK}^g)'$ is the random vector representing sequence i of group g . Consider \mathbf{X}_i^g and $\mathbf{X}_j^{g'}$.

Definition

The *Hamming Distance* $D_{ij}^{(g,g')}$ between sequences i of group g and j of group g' is

$$\begin{aligned} D_{ij}^{(g,g')} &= \frac{1}{K} \sum_{k=1}^K I(X_{ik}^g \neq X_{jk}^{g'}) \\ &= \frac{1}{K} \times \text{number of positions at which } \mathbf{X}_i^g \text{ and } \mathbf{X}_j^{g'} \text{ differ} \end{aligned}$$

where I is the indicator function (i.e., $I(A) = 1$ when A is true, and 0 otherwise), and when $g = g'$,

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

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

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

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

$$\bar{D}^{(g,g')} = \frac{1}{N_g N_{g'}} \sum_{i=1}^{N_g} \sum_{j=1}^{N_{g'}} D_{ij}^{(g,g')} = \frac{1}{N_g N_{g'} K} \sum_{i=1}^{N_g} \sum_{j=1}^{N_{g'}} \sum_{k=1}^K I(X_{ik}^g \neq X_{jk}^{g'})$$

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

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

which is a linear combination of U-statistics.

The Total Sum of Squares can be decomposed as

$$TSS = \sum_{g=1}^G \sum_{1 \leq i < j \leq N_g} (D_{ij}^g - \bar{D}.)^2 + \sum_{1 \leq g < g' \leq G} \sum_{i=1}^{N_g} \sum_{j=1}^{N_{g'}} (D_{ij}^{(g,g')} - \bar{D}.)^2 \quad (2.1)$$

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

where WSS , BSS , $AWSS$ and $ABSS$ stand, respectively, for the within-, between-, across/within- and across/between-group sum of squares.

3. Connections Between Sums of Squares and U-statistics

There are G groups with N_g sequences each. Yet, we can disregard the group clustering and think of the sequences as a random sample of size $\sum_{g=1}^G N_g = M$. Then

$$\begin{aligned} TSS &= \sum_{1 \leq i < j \leq M} (D_{ij} - \bar{D}.)^2 \\ &= \left(\frac{M(M-1)}{2} - 1 \right) \left(\frac{M(M-1)}{2} \right)^{-1} \frac{1}{2} \sum_{\substack{i < j, i' < j' \\ i \leq i' \text{ or } j \leq j'}} (D_{ij} - D_{i'j'})^2 \end{aligned} \quad (3.1)$$

$$\begin{aligned}
WSS &= \sum_{g=1}^G \sum_{1 \leq i < j \leq N_g} (D_{ij}^g - \bar{D}^g)^2 \\
&= \sum_{g=1}^G \left(\frac{N_g(N_g-1)}{2} - 1 \right) \left(\frac{N_g(N_g-1)}{2} \right)^{-1} \frac{1}{2} \sum_{\substack{i < j, i' < j' \\ i \leq i' \text{ or } j \leq j'}} (D_{ij}^g - D_{i'j'}^g)^2
\end{aligned} \tag{3.2}$$

and

$$\begin{aligned}
AWSS &= \sum_{1 \leq g < g' \leq G} \sum_{i=1}^{N_g} \sum_{j=1}^{N_{g'}} (D_{ij}^{(g,g')} - \bar{D}^{(g,g')})^2 \\
&= \sum_{1 \leq g < g' \leq G} (N_g N_{g'} - 1) \binom{N_g N_{g'}}{2}^{-1} \frac{1}{2} \sum_{i=1}^{N_g} \sum_{j=1}^{N_{g'}} \sum_{i'=1}^{N_g} \sum_{j'=1}^{N_{g'}} \sum_{\substack{i \leq i' \text{ or } j \leq j'}} (D_{ij}^{(g,g')} - D_{i'j'}^{(g,g')})^2
\end{aligned} \tag{3.3}$$

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

4. Asymptotic Distributions and Decompositions of U-statistics

Let F denote the distribution function of X_i and U^m be a U-statistic of degree m , computed from a sample of size n , with kernel $\phi(X_1, \dots, X_m)$ and $E(U^m) = \theta(F) = \theta$.

$$U^m \equiv U(X_1, \dots, X_n) = \binom{n}{m}^{-1} \sum_{1 \leq i_1 < \dots < i_m \leq n} \phi(X_{i_1}, \dots, X_{i_m}), \quad n \geq m \tag{4.1}$$

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

Let

$$\Psi_c(x_1, \dots, x_c) \equiv E\{\phi(x_1, \dots, x_c, X_{c+1}, \dots, X_m)\} \tag{4.2}$$

$$\psi_c(x_1, \dots, x_c) \equiv E\{\phi(x_1, \dots, x_c, X_{c+1}, \dots, X_m) - \theta\} \tag{4.3}$$

$$\xi_c \equiv E\{\Psi_c^2(X_1, \dots, X_c)\} - \theta^2 \quad \text{and} \quad \xi_0 \equiv 0. \tag{4.4}$$

The function Ψ_c has the following properties (Lee, 1990, p. 11):

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

Now

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

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 (4.4), each term in $\sum^{(c)}$ is equal to ξ_c . The number of terms in $\sum^{(c)}$ is

$$\frac{n(n-1)\cdots(n-2m+c+1)}{c!(m-c)!(m-c)!} = \binom{m}{c} \binom{n-m}{m-c} \binom{n}{m} \quad (4.5)$$

Since $\xi_0 = 0$,

$$\text{Var}(U^m) = \binom{n}{m}^{-1} \sum_{c=1}^m \binom{m}{c} \binom{n-m}{m-c} \xi_c \quad (4.6)$$

The inequality $0 \leq \xi_c \leq \frac{c}{d} \xi_d$ with $1 \leq c < d \leq m$ (Hoeffding, 1948) leads to

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

From (4.6) and (4.5), $n\text{Var}(U^m)$ is a decreasing function of n which tends to its lower bound $m^2\xi_1$ as n increases, i.e.,

$$\text{Var}(U^m) = \frac{m^2}{n} \xi_1 + O(n^{-2}) \quad (4.7)$$

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

$$n^{1/2}(U^m - \theta) \xrightarrow{d} N(0, m^2\xi_1) \quad (4.8)$$

(Hoeffding, 1948).

Definition 1

$F_n(x)$ is the empirical distribution function (d.f.)

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \epsilon(x - X_i) \quad x \in \mathbb{R}^p, \quad n \geq 1$$

with $\epsilon(u)$ being 1 if all p coordinates of u are nonnegative and 0 otherwise. ■

We may rewrite (4.1) as

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

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

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

$$U^m = \theta(F) + \sum_{h=1}^m \binom{m}{h} U_h^m \quad n \geq m \quad (4.9)$$

where $U_h^m = n^{-[h]} \sum_{1 \leq i_1 \neq \dots \neq i_h \leq n} \int_{\mathbb{R}^{ph}} \cdots \int \Psi_h(x_1, \dots, x_h) \prod_{j=1}^h d[\epsilon(x_j - X_{i_j}) - F(x_j)]$ for $1 \leq h \leq m$. Further, if we write

$$\begin{aligned} \Psi_h^\circ(x_1, \dots, x_h) &= \Psi_h(x_1, \dots, x_h) - \sum_{j=1}^h \Psi_{h-1}(x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_h) \\ &\quad + \cdots + (-1)^h \theta(F), \quad \forall (x_1, \dots, x_h) \in \mathbb{R}^{ph}, \end{aligned} \quad (4.10)$$

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

$$U_h^m = \binom{n}{h}^{-1} \sum_{1 \leq i_1 < \dots < i_h \leq n} \Psi_h^\circ(X_{i_1}, \dots, X_{i_h}), \quad 1 \leq h \leq m \quad (4.11)$$

and the U_h^m are themselves U-statistics. Note that for $h = 2$, we have

$$\begin{aligned} \mathbb{E}(U_2^m) &= \mathbb{E}(\Psi_2^\circ(X_1, X_2)) = \mathbb{E}(\Psi_2(X_1, X_2)) - \mathbb{E}(\Psi_1(X_1)) \\ &\quad - \mathbb{E}(\Psi_1(X_2)) + \theta(F) \\ &= \theta(F) - \theta(F) - \theta(F) + \theta(F) = 0 \end{aligned}$$

Let

$$\Psi_{h,h-1}^\circ(x_1, \dots, x_{h-1}) = \mathbb{E}[\Psi_h^\circ(X_1, \dots, X_{h-1}, X_h) \mid X_1, \dots, X_{h-1}]$$

Then

$$\begin{aligned} \Psi_{21}^\circ(X_1) &\equiv \mathbb{E}[\Psi_2^\circ(X_1, X_2) \mid X_1] \\ &= \mathbb{E}[\Psi_2(X_1, X_2) \mid X_1] - \mathbb{E}[\Psi_1(X_1) \mid X_1] - \mathbb{E}[\Psi_1(X_2) \mid X_1] + \theta(F) \\ &= \Psi_1(X_1) - \Psi_1(X_1) - \mathbb{E}(\Psi_1(X_2)) + \theta(F) \\ &= \theta(F) - \theta(F) = 0 \end{aligned}$$

$\xi_1^\circ \equiv \mathbb{E}[\Psi_{21}^\circ(X_1)]^2 - (\mathbb{E}(U_2^m))^2 = 0$ and by (4.6),

$$\begin{aligned} \text{Var}(U_2^m) &= \frac{4(n-2)}{n(n-1)} \xi_1^\circ + \frac{2\xi_2^\circ}{n(n-1)} \\ &= \frac{4}{n} \xi_1^\circ + O(n^{-2}) = O(n^{-2}) \end{aligned} \quad (4.12)$$

Consequently, $U_2^m = O_p(n^{-1})$.

From direct computation, $\mathbb{E}(U_h^m) = 0$, $\forall 1 \leq h \leq m$ and

$$\text{Var}(U_h^m) = \mathbb{E}[(U_h^m)^2] = O(n^{-h}), \quad h = 1, 2, \dots, m; \quad (4.13)$$

and we can write

$$U^m = \theta(F) + \frac{m}{n} \sum_{i=1}^n [\Psi_1(X_i) - \theta(F)] + O_p(n^{-1}) \quad (4.14)$$

Now we consider multiple-sample U-statistics.

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 $F^{(j)}(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^p$, 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

$$\theta(\mathbf{F}) = \int_{\mathbb{R}^{m_0}} \dots \int \phi(\mathbf{x}_i^{(j)}, 1 \leq i \leq m_j, 1 \leq j \leq c) \prod_{j=1}^c \prod_{i=1}^{m_j} dF^{(j)}(\mathbf{x}_i^{(j)}) \quad (4.15)$$

Definition 2

For a set of samples of sizes $\mathbf{n} = (n_1, n_2, \dots, n_c)$ with $n_j \geq m_j$, $1 \leq j \leq c$, $\mathbf{m} = (m_1, m_2, \dots, m_c)$ the *generalized U-statistic* for $\theta(\mathbf{F})$ is

$$U^{(\mathbf{m})} = \prod_{j=1}^c \binom{n_j}{m_j}^{-1} \sum_{(\mathbf{n})}^* \phi(\mathbf{X}_\alpha^{(j)}, \alpha = i_{j1}, \dots, i_{jm_j}, 1 \leq j \leq c), \quad (4.16)$$

where the summation $\sum_{(\mathbf{n})}^*$ extends over all $1 \leq i_{j1} < \dots < i_{jm_j} \leq n_j$, $1 \leq j \leq c$. $U^{(\mathbf{m})}$ 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

$$\Psi_{d_1 \dots d_c}(\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_{d_j}^{(j)}, 1 \leq j \leq c) \equiv E(\phi(\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_{d_j}^{(j)}, \mathbf{X}_{d_j+1}^{(j)}, \dots, \mathbf{X}_{m_j}^{(j)}, 1 \leq j \leq c)) \quad (4.17)$$

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

$$\xi_{\mathbf{d}}(\mathbf{F}) = E\left(\Psi_{\mathbf{d}}^2(\mathbf{X}_1^{(j)}, \dots, \mathbf{X}_{d_j}^{(j)}, 1 \leq j \leq c)\right) - \theta^2(\mathbf{F}), \quad \mathbf{0} \leq \mathbf{d} \leq \mathbf{m} \quad (4.18)$$

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

$$\text{Var} [U^{(\mathbf{m})}] = \sum_{j=1}^c n_j^{-1} \sigma_j^2 [1 + O(n_0^{-1})] \quad (4.19)$$

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

$$\sigma_j^2 = m_j^2 \xi_{\delta_{j1}, \dots, \delta_{jc}}(\mathbf{F}) \quad j = 1, \dots, c \quad (4.20)$$

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

For a two-sample U-statistic of degree (m_1, m_2) , the kernel is

$$\phi(X_1, \dots, X_{m_1}; Y_1, \dots, Y_{m_2})$$

and

$$U^{(m_1, m_2)} = \binom{n_1}{m_1}^{-1} \binom{n_2}{m_2}^{-1} \sum_{(\mathbf{n}, \mathbf{m})}^* \phi(X_{i_{11}}, \dots, X_{i_{1m_1}}; Y_{i_{21}}^{(2)}, \dots, Y_{i_{2m_2}})$$

where $\sum_{(\mathbf{n}, \mathbf{m})}^*$ extends over all $1 \leq i_{j1} < \dots < i_{jm_j} \leq n_j$, $j = 1, 2$. $U^{(m_1, m_2)}$ is an unbiased estimator of $\theta(F^{(1)}, F^{(2)})$.

Define

$$\begin{aligned} & \theta(F_{n_1}^{(1)}, F_{n_2}^{(2)}) \\ &= n_1^{-m_1} n_2^{-m_2} \times \sum_{i_{11}=1}^{n_1} \dots \sum_{i_{1m_1}=1}^{n_1} \sum_{i_{21}=1}^{n_2} \dots \sum_{i_{2m_2}=1}^{n_2} \phi(X_{i_{11}}, \dots, X_{i_{1m_1}}, Y_{i_{21}}, \dots, Y_{i_{2m_2}}) \end{aligned}$$

Then,

$$|U^{(m_1, m_2)} - \theta(F_{n_1}^{(1)}, F_{n_2}^{(2)})| = O_p(n_0^{-1}); \quad n_0 = \min(n_1, n_2) \quad (4.21)$$

provided the variance of $U^{(m_1, m_2)}$ exists.

$$\begin{aligned} & \Psi_{d_1 d_2}(x_1, \dots, x_{d_1}, y_1, \dots, y_{d_2}) \\ &= E\{\phi(x_1, \dots, x_{d_1}, X_{d_1+1}, \dots, X_{m_1}; y_1, \dots, y_{d_2}, Y_{d_2+1}, \dots, Y_{m_2})\}, \\ & \xi_{d_1 d_2} = E\{\Psi_{d_1 d_2}^2(X_1, \dots, X_{d_1}, Y_1, \dots, Y_{d_2})\} - \theta^2(F^{(1)}, F^{(2)}) \end{aligned} \quad (4.22)$$

for $d_1 = 0, \dots, m_1$, $d_2 = 0, \dots, m_2$. ($\xi_{00} \equiv 0$). Then,

$$(\gamma_{n_1, n_2})^{-1}(U^{(m_1, m_2)} - \theta(F^{(1)}, F^{(2)})) \xrightarrow{d} N(0, 1), \quad (4.23)$$

where

$$\gamma_{n_1, n_2}^2 = \left(\frac{m_1^2}{n_1}\right) \xi_{10} + \left(\frac{m_2^2}{n_2}\right) \xi_{01}. \quad (4.24)$$

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

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

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

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

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

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

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

5. Combining the U-statistics

We can write

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

where

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

$$\mathbf{U}_{1,3}^{(3)} = \binom{N_g}{3}^{-1} \sum_{i < j < j'} (D_{ij}^g - D_{jj'}^g)^2$$

are one-sample U-statistics of degree 3 and

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

$$\mathbf{U}_{2,3}^{(4)} = \binom{N_g}{4}^{-1} \sum_{i < i' < j' < j} (D_{ij}^g - D_{i'j'}^g)^2$$

are one-sample U-statistics of degree 4.

Note that $U_{j,k}^{(m)}$ represents a U-statistic of degree m with expected value μ_j , i.e., $E(U_{j,k}^{(m)}) = \mu_j$ for any k . The expected value of WSS is

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

where $\mu_{1g} = E(\mathbf{U}_{1,1}^{(3)}) = E(\mathbf{U}_{1,2}^{(3)}) = E(\mathbf{U}_{1,3}^{(3)}) = E(D_{ij}^g - D_{ij'}^g)^2$ and $\mu_{2g} = E(\mathbf{U}_{2,1}^{(4)}) = E(\mathbf{U}_{2,2}^{(4)}) = E(\mathbf{U}_{2,3}^{(4)}) = E[(D_{ij}^g - D_{i'j'}^g)^2]$.

Under H_0 , there is homogeneity among groups, i.e., for any g , $\theta_k^g = \theta_k$ and $\theta_{k_1 k_2}^g = \theta_{k_1 k_2}$. Therefore, $\mu_{1g} = \mu_1$ and $\mu_{2g} = \mu_2$, thus

$$E_0(WSS) = \sum_{g=1}^G (N_g - 2) \left\{ \mu_1 + \frac{(N_g - 3)}{4} \mu_2 \right\}$$

where

$$\mu_{1g} = \mu_1 = \frac{2}{K^2} \left[\sum_{k=1}^K \theta_k + \sum_{k_1 \neq k_2} \theta_{k_1 k_2} - \sum_{k=1}^K \theta_k(i, j; i, j') - \sum_{k_1 \neq k_2} \theta_{k_1 k_2}(i, j; i, j') \right] \quad (5.1)$$

and

$$\mu_{2g} = \mu_2 = \frac{2}{K^2} \left\{ \sum_{k=1}^K \theta_k(1 - \theta_k) + \sum_{k_1 \neq k_2} (\theta_{k_1 k_2} - \theta_{k_1} \theta_{k_2}) \right\} \quad (5.2)$$

with

$$\theta_k = \theta_k(i, j) = P(X_{ik} \neq X_{jk}) = \sum_{c=0}^{C-1} p_k(c) [1 - p_k(c)] \quad (5.3)$$

$$\begin{aligned} \theta_{k_1 k_2} &= \theta_{k_1, k_2}(i, j) = P(X_{ik_1} \neq X_{jk_1}; X_{ik_2} \neq X_{jk_2}) \\ &= \sum_{c_1, c_2=0}^{C-1} p_{k_1 k_2}(c_1, c_2) \left[\sum_{\substack{c_3=0 \\ c_3 \neq c_1}}^{C-1} \sum_{\substack{c_4=0 \\ c_4 \neq c_2}}^{C-1} p_{k_1 k_2}(c_3, c_4) \right] \end{aligned} \quad (5.4)$$

$$\theta_k^g(i, j; i, j') = P(X_{ik}^g \neq X_{jk}^g, X_{ik}^g \neq X_{j'k}^g) = \sum_{c=0}^{C-1} p_k^g(c)[1 - p_k^g(c)]^2, \quad (5.5)$$

$$\begin{aligned} \theta_{k_1 k_2}^g(i, j; i, j') &= P(X_{ik_1}^g \neq X_{jk_1}^g, X_{ik_2}^g \neq X_{j'k_2}^g) \\ &= \sum_{c_1, c_2=0}^{C-1} p_{k_1 k_2}^g(c_1, c_2)[1 - p_{k_1}^g(c_1)][1 - p_{k_2}^g(c_2)], \end{aligned} \quad (5.6)$$

$$p_k(c) = P(X_{ik}^g = c) \quad \text{and} \quad p_{k_1 k_2}^g(c_1, c_2) = P(X_{ik_1}^g = c_1, X_{ik_2}^g = c_2) \quad (5.7)$$

Decomposing WSS , under H_0 ,

$$\begin{aligned} WSS &= \sum_{g=1}^G (N_g - 2) \left(\mu_1 + \frac{(N_g - 3)}{4} \mu_2 \right) \\ &\quad + \sum_{g=1}^g (N_g - 2) \frac{3}{N_g} \sum_{i=1}^{N_g} [\Psi_{(1)1}(\mathbf{X}_i) - \mu_1] + O_p(1) \\ &\quad + \sum_{g=1}^G \frac{(N_g - 2)(N_g - 3)}{N_g} \sum_{i=1}^{N_g} [\Psi_{(2)1}(\mathbf{X}_i) - \mu_2] + O_p(N_g) \end{aligned} \quad (5.8)$$

and the associated mean square expression is

$$\begin{aligned} WMS &\equiv \frac{WSS}{\sum_{g=1}^G \binom{N_g}{2}} = \frac{2WSS}{\sum_{g=1}^G N_g(N_g - 1)} \\ &= \frac{2}{\sum_{g=1}^G N_g(N_g - 1)} \left\{ \sum_{g=1}^G (N_g - 2) \left(\mu_1 + \frac{(N_g - 3)}{4} \mu_2 \right) \right. \\ &\quad + \sum_{g=1}^g (N_g - 2) \frac{3}{N_g} \sum_{i=1}^{N_g} [\Psi_{(1)1}(\mathbf{X}_i) - \mu_1] + O_p(1) \\ &\quad \left. + \sum_{g=1}^G \frac{(N_g - 2)(N_g - 3)}{N_g} \sum_{i=1}^{N_g} [\Psi_{(2)1}(\mathbf{X}_i) - \mu_2] + O_p(N_g) \right\} \\ &= \frac{\sum_{g=1}^G (N_g - 2)(N_g - 3)}{2 \sum_{g=1}^G N_g(N_g - 1)} \left\{ \mu_2 + \frac{4}{N_g} \sum_{i=1}^{N_g} (\Psi_{(2)1}(\mathbf{X}_i) - \mu_2) \right\} + O_p(N_0^{-1}) \end{aligned} \quad (5.9)$$

with

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

and

$$\begin{aligned} \text{Var}_0(WMS) &= \frac{4 \sum_{g=1}^G [N_g^3 - 10N_g^2 + 37N_g - 60 + 36/N_g]}{[\sum_{g=1}^G N_g(N_g - 1)]^2} \xi_1^{(2)} + O(N_g^{-2}) \\ &= \frac{4 \sum_{g=1}^G N_g^3}{[\sum_{g=1}^G N_g(N_g - 1)]^2} \xi_1^{(2)} + O(N_g^{-2}) \end{aligned} \quad (5.10)$$

For *AWSS*,

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

where

$$\mathbf{U}_{4,1}^{(2,2)} = \left[\binom{N_g}{2} \binom{N_{g'}}{2} \right]^{-1} \sum_{\substack{i \neq i' \\ j \neq j'}} \sum_{\substack{j \neq j' \\ i \neq i'}} (D_{ij}^{(g,g')} - D_{i'j'}^{(g,g')})^2 \quad \text{and}$$

$$\mathbf{U}_{4,2}^{(2,2)} = \left[\binom{N_g}{2} \binom{N_{g'}}{2} \right]^{-1} \sum_{\substack{i \neq j \\ i \neq i'}} \sum_{\substack{j \neq j' \\ j \neq i'}} (D_{ij}^{(g,g')} - D_{i'j'}^{(g,g')})^2$$

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

$$\mathbf{U}_{5,1}^{(1,2)} = \left[\binom{N_g}{1} \binom{N_{g'}}{2} \right]^{-1} \sum_{\substack{i=1 \\ i \neq i'}}^{N_g} \sum_{\substack{1 \leq j, j' \leq N_{g'} \\ j \neq j'}} (D_{ij}^{(g,g')} - D_{i'j'}^{(g,g')})^2 \quad \text{and}$$

$$\mathbf{U}_{5,2}^{(2,1)} = \left[\binom{N_g}{2} \binom{N_{g'}}{1} \right]^{-1} \sum_{\substack{j=1 \\ j \neq i'}}^{N_{g'}} \sum_{\substack{1 \leq i, i' \leq N_g \\ i \neq i'}} (D_{ij}^{(g,g')} - D_{i'j}^{(g,g')})^2$$

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

$$E(AWSS) = \sum_{1 \leq g < g' \leq G} \left[\frac{(N_g - 1)(N_{g'} - 1)}{2} \mu_{4(g,g')} + \frac{(N_g - 1)}{2} \mu_{5(g,g')} + \frac{(N_{g'} - 1)}{2} \mu_{5(g,g')} \right]$$

with $\mu_{4(g,g')} = E(\mathbf{U}_{4,1}^{(2,2)}) = E(\mathbf{U}_{4,2}^{(2,2)})$ and $\mu_{5(g,g')} = E(\mathbf{U}_{5,1}^{(1,2)}) = E(\mathbf{U}_{5,2}^{(2,1)})$.

Under H_0 , $\mu_{4(g,g')} = \mu_4$ and $\mu_{5(g,g')} = \mu_5$, therefore

$$E_0(AWSS) = \frac{1}{2} \sum_{1 \leq g < g' \leq G} [(N_g - 1)(N_{g'} - 1)\mu_4 + (N_g + N_{g'} - 2)\mu_5]$$

where $\mu_4 = \mu_2$ is given by (5.2) and $\mu_5 = \mu_1$ is given by (5.1).

AWSS can be decomposed as

$$\begin{aligned} AWSS &= \sum_{1 \leq g < g' \leq G} \left[\frac{(N_g - 1)(N_{g'} - 1)}{2} \left(\mu_{4(g,g')} + \frac{2}{N_g} \sum_{i=1}^{N_g} (\Psi_{(4)10}(\mathbf{X}_i^g) - \mu_{4(g,g')}) \right) \right. \\ &\quad \left. + \frac{2}{N_{g'}} \sum_{j=1}^{N_{g'}} (\Psi_{(4)01}(\mathbf{X}_j^{g'}) - \mu_{(g,g')4}) + O_p(N_0^{-1}) \right) \\ &\quad + \frac{(N_g - 1)}{2} \left(\mu_{5(g,g')} + \frac{2}{N_g} \sum_{i=1}^{N_g} (\Psi_{(5)10}(\mathbf{X}_i^g) - \mu_{5(g,g')}) \right) \\ &\quad \left. + \frac{1}{N_{g'}} \sum_{j=1}^{N_{g'}} (\Psi_{(5)01}(\mathbf{X}_j^{g'}) - \mu_{5(g,g')}) \right) \end{aligned}$$

$$\begin{aligned}
& + \frac{(N_{g'} - 1)}{2} \left(\mu_{5(g,g')} + \frac{1}{N_g} \sum_{i=1}^{N_g} (\Psi_{(5)10}(\mathbf{X}_i^g) - \mu_{5(g,g')}) \right. \\
& \left. + \frac{2}{N_{g'}} \sum_{j=1}^{N_{g'}} (\Psi_{(5)01}(\mathbf{X}_j^{g'}) - \mu_{5(g,g')}) + O_p(N_0^{-1}) \right) \Big] \tag{5.11}
\end{aligned}$$

with $N_0 = \min_{1 \leq g \leq G} (N_g)$.

The associated mean-square expression is

$$\begin{aligned}
AWMS &= \frac{AWSS}{\sum_{g < g'} N_g N_{g'}} \\
&= \frac{1}{2} \sum_{1 \leq g < g' \leq G} \frac{(N_g - 1)(N_{g'} - 1)}{\left[\sum_{g < g'} N_g N_{g'} \right]} \left[\mu_{4(g,g')} + \frac{2}{N_g} \sum_{i=1}^{N_g} (\Psi_{(4)10}(\mathbf{X}_i^g) - \mu_{4(g,g')}) \right. \\
& \left. + \frac{2}{N_{g'}} \sum_{j=1}^{N_{g'}} (\Psi_{(4)01}(\mathbf{X}_j^{g'}) - \mu_{4(g,g')}) \right] + O_p(N_0^{-1}) \tag{5.12}
\end{aligned}$$

$$E_0(AWMS) = \frac{\mu_4}{2} + O(N_0^{-1})$$

Note that $E_0(AWMS) = E_0(WMS)$, since under H_0 , $\mu_4 = \mu_2$.

$$\begin{aligned}
\text{Var}_0(AWMS) &= \frac{1}{(\sum_{g < g'} N_g N_{g'})^2} \left\{ \sum_{g < g'} \frac{(N_g - 1)^2 (N_{g'} - 1)^2}{4} \left(\frac{4}{N_g} \xi_{10}^{(4)} + \frac{4}{N_{g'}} \xi_{01}^{(4)} \right) \right\} + O(N_0^{-2}) \\
&= \frac{1}{(\sum_{g < g'} N_g N_{g'})^2} \left\{ \sum_{g < g'} (N_g - 1)^2 (N_{g'} - 1)^2 \left(\frac{\xi_{10}^{(4)}}{N_g} + \frac{\xi_{01}^{(4)}}{N_{g'}} \right) \right\} + O(N_0^{-2}) \tag{5.13}
\end{aligned}$$

Now

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

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

$$\mathbf{D}_1 = \left(\sqrt{N_1(N_1 - 1)}(\bar{D}^1 - \bar{D}.) \dots \sqrt{N_G(N_G - 1)}(\bar{D}^G - \bar{D}.) \right)'$$

Note that

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

$$E(\bar{D}^{(g,g')}) = \frac{1}{K} \sum_{k=1}^K \theta_k^{(g,g')} = \bar{\theta}^{(g,g')}$$

and

$$\mathbb{E}(\bar{D}.) = \binom{M}{2}^{-1} \left[\sum_{g=1}^G \frac{N_g(N_g-1)}{2} \bar{\theta}^g + \sum_{g < g'} N_g N_{g'} \bar{\theta}^{(g,g')} \right]$$

where $M = \sum_{g=1}^G N_g$.

Therefore,

$$\nu_1 \equiv \mathbb{E}(\bar{D}^g - \bar{D}.) = \bar{\theta}^g - \binom{M}{2}^{-1} \left[\sum_{g=1}^G \frac{N_g(N_g-1)}{2} \bar{\theta}^g + \sum_{g < g'} N_g N_{g'} \bar{\theta}^{(g,g')} \right] \quad (5.14)$$

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

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

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

$$\text{Var}(\bar{D}^{(g,g')}) = \frac{1}{N_g} \xi_{10}^{(13)} + \frac{1}{N_{g'}} \xi_{01}^{(13)} + O(N_0^{-2}) \quad (5.15)$$

where $\xi_{10}^{(13)} \equiv \mathbb{E}[\psi_{(13)10}^2(\mathbf{X}_i^g)]$ and $\xi_{01}^{(13)} \equiv \mathbb{E}[\psi_{(13)01}^2(\mathbf{X}_j^{g'})]$. 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 $\theta_k^{(g,g')} = \theta_k^g = \theta_k$. Therefore, under H_0 ,

$$\sqrt{N_g} (\bar{D}^g - \bar{\theta}.) \xrightarrow{d} \mathbb{N}(0, 4\xi_1^{(12)}) \quad (5.16)$$

and

$$\gamma_{13}^{-1} (\bar{D}^{(g,g')} - \bar{\theta}.) \xrightarrow{d} \mathbb{N}(0, 1) \quad (5.17)$$

where $\gamma_{13}^2 = \frac{1}{N_g} \xi_{10}^{(13)} + \frac{1}{N_{g'}} \xi_{01}^{(13)} = \left(\frac{1}{N_g} + \frac{1}{N_{g'}} \right) \xi_1^{(12)}$ by (5.15) and under H_0 .

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

$$\bar{D} = \binom{M}{2}^{-1} \left[\sum_{g=1}^G \frac{N_g(N_g-1)}{2} \bar{D}^g + \sum_{1 \leq g < g' \leq G} N_g N_{g'} \bar{D}^{(g,g')} \right]$$

where $M = \sum_{g=1}^G N_g$. Under H_0 ,

$$\eta_1 \equiv \mathbb{E}_0(\bar{D}.) = \binom{M}{2}^{-1} \left[\sum_{g=1}^G \binom{N_g}{2} + \sum_{1 \leq g < g' \leq G} N_g N_{g'} \right] \bar{\theta} = \bar{\theta}.$$

$$\begin{aligned}
\sigma_1^2 &\equiv \text{Var}_0(\bar{D}.) \\
&= \frac{4}{M^2(M-1)^2} \left\{ \sum_{g=1}^G N_g(N_g-1)^2 \xi_1^{(12)} + \sum_{1 \leq g < g' \leq G} N_g^2 N_{g'}^2 \left(\frac{1}{N_g} \xi_{10}^{(13)} + \frac{1}{N_{g'}} \xi_{01}^{(13)} \right) \right. \\
&\quad \left. + 2 \sum_{\substack{1 \leq g, g', g^\dagger \leq G \\ g \neq g' \neq g^\dagger}} N_g N_{g'} N_{g^\dagger} \xi_{10}^{(13,1;13,2)} + 2 \sum_{g=1}^G \sum_{g'=g+1}^G N_g(N_g-1) N_{g'} \xi_1^{(12,13)} \right\}
\end{aligned}$$

where $\xi_{10}^{(13,1;13,2)} = \text{E}\{\psi_{(13,1)10}(\mathbf{X}_i^g) \psi_{(13,2)10}(\mathbf{X}_i^g)\}$ and $\psi_{(13,2)10}(\mathbf{x}_i^g) = \text{E}[\phi_{13,2}(\mathbf{x}_i^g, \mathbf{X}_j^{g^\dagger}) - \bar{\theta}^{(g,g^\dagger)}]$. Under H_0 , $\psi_{(12)1}(\mathbf{X}_i) = \psi_{(13)10}(\mathbf{X}_i) = \psi_{(13)01}(\mathbf{X}_i) = \psi_{(13,1)10}(\mathbf{X}_i) = \psi_{(13,2)10}(\mathbf{X}_i)$. Therefore, $\xi_1^{(12)} = \xi_{10}^{(13)} = \xi_{01}^{(13)} = \xi_{10}^{(13,1;13,2)} = \xi_1^{(12,13)}$ and

$$\begin{aligned}
\sigma_1^2 &= \frac{4\xi_1^{(12)}}{M^2(M-1)^2} \left\{ \sum_{g=1}^G N_g(N_g-1)^2 + \sum_{1 \leq g < g' \leq G} N_g N_{g'} (N_g + N_{g'}) \right. \\
&\quad \left. + 2 \sum_{\substack{1 \leq g, g', g^\dagger \leq G \\ g \neq g' \neq g^\dagger}} N_g N_{g'} N_{g^\dagger} + 2 \sum_{g=1}^G \sum_{g'=g+1}^G N_g(N_g-1) N_{g'} \right\} \quad (5.18)
\end{aligned}$$

Hence, under H_0 ,

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

Now, under H_0 ,

$$\nu_1 = \text{E}_0(\bar{D}^g - \bar{D}.) = \bar{\theta}. - \bar{\theta}. = 0 \quad (5.19)$$

and

$$\begin{aligned}
\tau_{1g}^2 &\equiv \text{Var}_0(\bar{D}^g - \bar{D}.) \\
&= \text{Var}_0(\bar{D}^g) + \text{Var}(\bar{D}.) - 2\text{Cov}_0(\bar{D}^g, \bar{D}.) \\
&= \left[\frac{4}{N_g} - \frac{8(N_g-1)}{M(M-1)} \right] \xi_1^{(12)} + \sigma_1^2 - \frac{8(M-N_g)}{M(M-1)} \xi_1^{(12,13)} \quad (5.20)
\end{aligned}$$

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

Then,

$$\begin{aligned}
\tau_{1g}^2 &= 4\xi_1^{(12)} \left\{ \frac{1}{N_g} + \frac{\sum_{g=1}^G N_g(N_g-1)^2}{M^2(M-1)^2} + \sum_{g < g'} \frac{N_g N_{g'} (N_{g'} + N_g)}{M^2(M-1)^2} + 2 \sum_{g \neq g' \neq g^\dagger} \frac{N_g N_{g'} N_{g^\dagger}}{M^2(M-1)^2} \right. \\
&\quad \left. + 2 \sum_{g=1}^G \sum_{g'=g+1}^G \frac{N_g(N_g-1) N_{g'}}{M^2(M-1)^2} - 2 \frac{(N_g-1)}{M(M-1)} - 2 \frac{(M-N_g)}{M(M-1)} \right\} \quad (5.21)
\end{aligned}$$

So,

$$\tau_{1g}^{-1} (\bar{D}^g - \bar{D}.) \xrightarrow{d} \text{N}(0, 1)$$

Since BSS is a quadratic form of normal random variables,

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

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 $N_g(N_g - 1)\tau_{1g}^2$ and the off-diagonal elements, under H_0 , are

$$\begin{aligned} \sqrt{N_g(N_g - 1)N_{g'}(N_{g'} - 1)\tau_{1gg'}} &= \sqrt{N_g(N_g - 1)N_{g'}(N_{g'} - 1)\text{Cov}_0(\bar{D}^g - \bar{D}, \bar{D}^{g'} - \bar{D})} \\ &= \sqrt{N_g(N_g - 1)N_{g'}(N_{g'} - 1)} \left(-\frac{4\xi_1^{(12)}}{M(M-1)} [N_g - 1 + M - N_g + N_{g'} - 1 + M - N_{g'}] + \sigma_1^2 \right) \\ &= \sqrt{N_g(N_g - 1)N_{g'}(N_{g'} - 1)} \left(-\frac{8\xi_1^{(12)}}{M} + \sigma_1^2 \right) \end{aligned} \quad (5.22)$$

where σ_1^2 is given by (5.18).

Now,

$$E_0(BSS) = \frac{1}{2} \text{trace}(\boldsymbol{\Sigma}_1) = \frac{1}{2} \sum_{g=1}^G N_g(N_g - 1)\tau_{1g}^2$$

and

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

Let

$$BMS = \frac{BSS}{\sum_{g=1}^G \binom{N_g}{2}} = \frac{\mathbf{D}'_1 \mathbf{D}_1}{\sum_{g=1}^G N_g(N_g - 1)}$$

Then

$$E_0(BMS) = \frac{\text{trace}(\boldsymbol{\Sigma}_1)}{\sum_{g=1}^G N_g(N_g - 1)} = \frac{\sum_{g=1}^G N_g(N_g - 1)\tau_{1g}^2}{\sum_{g=1}^G N_g(N_g - 1)}$$

and

$$\text{Var}_0(BMS) = \frac{1}{\left[\sum_{g=1}^G \binom{N_g}{2} \right]^2} \text{Var}_0(BSS) = \frac{\text{trace}(\boldsymbol{\Sigma}_1)^2}{\left[\sum_{g=1}^G N_g(N_g - 1) \right]^2}$$

For $ABSS$ we have,

$$ABSS = \sum_{1 \leq g < g' \leq G} \sum_{i=1}^{N_g} \sum_{j=1}^{N_{g'}} (\bar{D}^{(g,g')} - \bar{D})^2 = \mathbf{D}'_2 \mathbf{D}_2$$

where $\mathbf{D}_2 = [\sqrt{N_1 N_2}(\bar{D}^{(1,2)} - \bar{D}), \sqrt{N_1 N_3}(\bar{D}^{(1,3)} - \bar{D}), \dots, \sqrt{N_{G-1} N_G}(\bar{D}^{(G-1,G)} - \bar{D})]'$ is a $\frac{G(G-1)}{2} \times 1$ vector.

Let

$$\nu_2 \equiv E(\bar{D}^{(g,g')} - \bar{D}) = \bar{\theta}^{(g,g')} - \binom{M}{2}^{-1} \left[\sum_{g=1}^G \frac{N_g(N_g - 1)}{2} \bar{\theta}^g + \sum_{g < g'} N_g N_{g'} \bar{\theta}^{(g,g')} \right]$$

Under H_0 ,

$$\nu_2 = E_0(\bar{D}^{(g,g')} - \bar{D}.) = \bar{\theta} - \bar{\theta} = 0 \quad (5.23)$$

and

$$\begin{aligned} \tau_{2(g,g')}^2 &\equiv \text{Var}(\bar{D}^{(g,g')} - \bar{D}.) \\ &= \text{Var}(\bar{D}^{(g,g')}) + \text{Var}(\bar{D}.) - 2\text{Cov}(\bar{D}^{(g,g')}, \bar{D}.) \\ &= \frac{1}{N_g} \xi_{10}^{(13)} + \frac{1}{N_{g'}} \xi_{01}^{(13)} + \sigma_1^2 - \frac{4}{M(M-1)} \left[(N_g - 1) \xi_1^{(12,13)} + (N_{g'} - 1) \xi_1^{(12,13)} + N_{g'} \xi_{01}^{(13)} \right. \\ &\quad \left. + 2(M - N_g - N_{g'}) \xi_{10}^{(13,1;13,2)} \right] \end{aligned} \quad (5.24)$$

Note that under H_0 there is homogeneity among groups,

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

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 ,

$$\tau_{2(g,g')}^2 = \left[\frac{(N_g + N_{g'})}{N_g N_{g'}} - \frac{8}{M} \right] \xi_1^{(12)} + \sigma_1^2 \quad (5.25)$$

As in BSS ,

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

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

$$\sqrt{N_g N_{g'} N_{g^\dagger} N_{g^\dagger}} \text{Cov}(\bar{D}^{(g,g')} - \bar{D}., \bar{D}^{(g^\dagger, g^\dagger)} - \bar{D}.) = \sqrt{N_g N_{g'} N_{g^\dagger} N_{g^\dagger}} \left(-\frac{8}{M} \xi_1^{(12)} + \sigma_1^2 \right) \quad (5.26)$$

and if there is one group in common, i.e., $g = g^\dagger$ or $g' = g^\dagger$,

$$N_g \sqrt{N_{g'} N_{g^\dagger}} \text{Cov}(\bar{D}^{(g,g')} - \bar{D}., \bar{D}^{(g,g^\dagger)} - \bar{D}.) = N_g \sqrt{N_{g'} N_{g^\dagger}} \left(\frac{\xi_1^{(12)}}{N_g} - \frac{8\xi_1^{(12)}}{M} + \sigma_1^2 \right) \quad (5.27)$$

Now

$$E_0(ABSS) = \text{trace}(\mathbf{\Sigma}_2) = \sum_{g < g'} N_g N_{g'} \tau_{2(g,g')}^2$$

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

The corresponding mean-square term is defined as

$$ABMS = \frac{ABSS}{\sum_{g < g'} N_g N_{g'}} = \frac{\mathbf{D}'_2 \mathbf{D}_2}{\sum_{g < g'} N_g N_{g'}}$$

Then

$$\mathbb{E}_0(ABMS) = \frac{\text{trace}(\boldsymbol{\Sigma}_2)}{\sum_{g < g'} N_g N_{g'}} = \frac{\sum_{g < g'} N_g N_{g'} \tau_{2(g, g')}^2}{\sum_{g < g'} N_g N_{g'}}$$

$$\text{Var}_0(ABMS) = \frac{\text{trace}(\boldsymbol{\Sigma}_2)^2}{[\sum_{g < g'} N_g N_{g'}]^2}$$

6. Test Statistics

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

$$\frac{WMS}{AWMS} = \frac{\frac{\sum_{g=1}^G (N_g - 2)(N_g - 3)}{2 \sum_{g=1}^G N_g (N_g - 1)} \left\{ \mu_2 + \frac{4}{N_g} \sum_{i=1}^{N_g} (\Psi_{(2)1}(\mathbf{X}_i) - \mu_2) \right\} + O_p(N_0^{-1})}{\frac{\sum_{g < g'} (N_g - 1)(N_{g'} - 1)}{2 \sum_{g < g'} N_g N_{g'}} \left\{ \mu_2 + \frac{2}{N_g} \sum_{i=1}^{N_g} (\Psi_{(2)1}(\mathbf{X}_i) - \mu_2) + \frac{2}{N_{g'}} \sum_{i=1}^{N_{g'}} (\Psi_{(2)1}(\mathbf{X}_i) - \mu_2) \right\} + O_p(N_0^{-1})}$$

Provided that $N_0 = \min_{1 \leq g \leq G} (N_g)$ and $N_g = O(N_0)$, $\forall g = 1, \dots, G$, we have that $\frac{WMS}{AWMS} \xrightarrow{p} 1$ as $N_0 \rightarrow \infty$, i.e., asymptotically the distribution of $\frac{WMS}{AWMS}$ is degenerate.

By (5.21) and (5.22) we have that $\boldsymbol{\Sigma}_1 = O(N_0)$ and by (5.25), (5.26) and (5.27), $\boldsymbol{\Sigma}_2 = O(N_0)$.

$$BMS = \frac{BSS}{\sum_{g=1}^G \binom{N_g}{2}} \sim \frac{\sum_{g=1}^G \lambda_{1g} (\chi_1^2)_g}{\sum_{g=1}^G N_g (N_g - 1)}$$

where λ_{1g} 's are the characteristic roots of $\text{Var}(\mathbf{D}_1) = \boldsymbol{\Sigma}_1$.

$$ABMS = \frac{ABSS}{\sum_{g < g'} N_g N_{g'}} \sim \frac{1}{\sum_{g < g'} N_g N_{g'}} \sum_{i=1}^{G(G-1)/2} \lambda_{2i} (\chi_1^2)_i$$

where λ_{2i} 's are the characteristic roots of $\text{Var}(\mathbf{D}_2) = \boldsymbol{\Sigma}_2$.

Also, under H_0 , by theoretical results pertaining to U-statistics

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

and

$$\sqrt{N_0}(AWMS - \mu_2/2) \rightarrow N\left(0, \xi_1^{(2)}\right) .$$

Thus,

$$BMS = O_p(N_0^{-1}) \quad \text{and} \quad ABMS = O_p(N_0^{-1})$$

while

$$WMS = O_p(N_g^{-1/2}) \quad \text{and} \quad AWMS = O_p(N_0^{-1/2})$$

Define

$$T_{N,2} \equiv N_0 \left(\frac{BMS}{WMS} \right) \quad \text{and} \quad T_{N,3} \equiv N_0 \left(\frac{ABMS}{AWMS} \right) .$$

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

$$T_{N,2} = \frac{2N_0(BMS)}{\mu_2} + O_p(N_0^{-1/2})$$

and

$$T_{N,3} = \frac{2N_0(ABMS)}{\mu_2} + O_p(N_0^{-1/2})$$

Therefore,

$$T_{N,2} \sim \frac{2N_0}{\mu_2} \frac{\sum_{g=1}^G \lambda_{1g} (\chi_1^2)_g}{\sum_{g=1}^G N_g(N_g - 1)}$$

and

$$T_{N,3} \sim \frac{2N_0}{\mu_2 \sum_{g < g'} N_g N_{g'}} \sum_{i=1}^{G(G-1)/2} \lambda_{2i} (\chi_1^2)_i$$

Because the elements of Σ_1 and Σ_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.

7. 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,

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

Let $\mathbf{V}_N = \mathbf{D}_1 - \nu_1$, where ν_1 is a vector $G \times 1$ with elements $\sqrt{N_g(N_g - 1)}\nu_1$, $g = 1, \dots, G$ and ν_1 is given by (5.14). Then, $E(\mathbf{V}_N) = \mathbf{0}$ and $\text{Var}(\mathbf{V}_N) = \text{Var}(\mathbf{D}_1) = \Sigma_1 = N_0 \Sigma_1^* = O(N_0)$. Therefore,

$$Q(\mathbf{D}_1) = \mathbf{D}'_1 \mathbf{D}_1 = \mathbf{V}'_N \mathbf{V}_N + 2\nu'_1 \mathbf{V}_N + \nu'_1 \nu_1$$

Since $\frac{\mathbf{V}_N}{\sqrt{N_0}} \sim N(\mathbf{0}, \Sigma_1^*)$,

$$\frac{\mathbf{V}'_N \mathbf{V}_N}{N_0} \sim \sum_{g=1}^G \lambda_g^* (\chi_1^2)_g$$

where λ_g^* are the characteristic roots of Σ_1^* . Also,

$$\frac{2\nu'_1 \mathbf{V}_N}{\sqrt{N_0}} \sim N(\mathbf{0}, 4\nu'_1 \Sigma_1^* \nu_1) \quad \text{or} \quad \frac{2\nu'_1 \mathbf{V}_N}{N_0^{3/2}} \sim N(\mathbf{0}, 4\nu_1^2 \Sigma_1^*)$$

With the above results, we can see that $\mathbf{V}_N = O_p(N_0^{1/2})$, $\mathbf{V}'_N \mathbf{V}_N = O_p(N_0)$ and $\nu'_1 \mathbf{V}_N = O_p(N_0^{3/2})$.

Now,

$$\begin{aligned} T_{N,2} &= \frac{2N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \mathbf{V}'_N \mathbf{V}_N + \frac{4N_0^{3/2} \nu'_1}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \left(\frac{\mathbf{V}_N}{\sqrt{N_0}} \right) + \frac{2N_0 \nu_1^2}{\mu_2} + O_p(N_0^{-1/2}) \\ \left(\frac{T_{N,2} - 2N_0 \nu_1^2 / \mu_2}{4N_0^{5/2} \nu_1 / [\mu_2 \sum_{g=1}^G N_g(N_g - 1)]} \right) &= \frac{\sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D}.) - \nu_1)^2}{2N_0^{3/2} \nu_1} \\ &+ \frac{1}{N_0^{3/2}} \sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D}.) - \nu_1 + O_p(N_0^{-1}) \end{aligned}$$

Note that

$$\frac{\sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D}.) - \nu_1)^2}{2N_0^{3/2} \nu_1} = O_p(N_0^{-1/2}),$$

since $\frac{\sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D}.) - \nu_1)^2}{N_0} = \frac{\mathbf{V}'_N \mathbf{V}_N}{N_0} = O_p(1)$

and

$$\frac{1}{N_0^{3/2}} \sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D}.) - \nu_1 = O_p(1),$$

since

$$\frac{\nu_1}{N_0^{3/2}} \sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D} - \nu_1) = \frac{\nu_1' \mathbf{V}_N}{N_0^{3/2}} = O_p(1)$$

So, for a fixed $\nu_1 \neq 0$, as $N_0 \rightarrow \infty$,

$$\left(\frac{T_{N,2} - 2N_0\nu_1^2/\mu_2}{4N_0^{5/2}\nu_1/[\mu_2 \sum_{g=1}^G N_g(N_g - 1)]} \right) = \frac{1}{N_0^{3/2}} \sum_{g=1}^G N_g(N_g - 1)(\bar{D}^g - \bar{D} - \nu_1) + O_p(N_0^{-1/2})$$

Thus,

$$P(T_{N,2} > \nu_1) = P\left(Z > \sum_{g=1}^G N_g(N_g - 1) \frac{(\mu_2 - 2N_0\nu_1)}{4N_0^{5/2}}\right) \rightarrow 1, \quad \text{as } N_0 \rightarrow \infty,$$

i.e., this test is consistent.

Now, consider a local alternative hypothesis. Let $\nu_1 = \frac{1}{\sqrt{N_0}}\gamma_1^*$, where γ_1^* is a constant. Then,

$$\begin{aligned} T_{N,2} &= \frac{2N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \mathbf{V}'_N \mathbf{V}_N \\ &+ \frac{4\gamma_1^*}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \left[\sqrt{N_0} \sum_{g=1}^G N_g(N_g - 1) \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N_0}}\gamma_1^* \right) \right] \\ &+ \frac{2}{\mu_2} (\gamma_1^*)^2 + O_p(N_0^{-1/2}) \end{aligned}$$

$$\begin{aligned} \left(\frac{T_{N,2} - 2(\gamma_1^*)^2/\mu_2}{4\gamma_1^* N_0^2 / \left(\mu_2 \sum_{g=1}^G N_g(N_g - 1) \right)} \right) &= \frac{\sum_{g=1}^G N_g(N_g - 1) \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N_0}}\gamma_1^* \right)^2}{2N_0\gamma_1^*} \\ &+ N_0^{-3/2} \sum_{g=1}^G N_g(N_g - 1) \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N_0}}\gamma_1^* \right) + O_p(N_0^{-1/2}) \end{aligned}$$

Note that

$$\frac{\sum_{g=1}^G N_g(N_g - 1) \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N_0}}\gamma_1^* \right)^2}{2N_0\gamma_1^*} = O_p(1) \quad \text{and} \quad \frac{\sum_{g=1}^G N_g(N_g - 1) \left(\bar{D}^g - \bar{D} - \frac{1}{\sqrt{N_0}}\gamma_1^* \right)}{N_0^{3/2}} = O_p(1)$$

Therefore, $T_{N,2}$ no longer follows a Normal distribution as $N_0 \rightarrow \infty$. It is a convolution of a linear combination of chi-square random variables and a normal random variable:

$$T_{N,2} = \frac{2N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \mathbf{V}'_N \mathbf{V}_N + \frac{4N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} (\gamma_1^*)' \mathbf{V}_N + \frac{2(\gamma_1^*)^2}{\mu_2} + O_p(N_0^{-1/2}),$$

where $\gamma_1^* = (\sqrt{N_1(N_1-1)}\frac{\gamma_1^*}{\sqrt{N_0}}, \dots, \sqrt{N_G(N_G-1)}\frac{\gamma_1^*}{\sqrt{N_0}})'$.

$$T_{N,2} \sim \frac{2N_0^2}{\mu_2 \sum_{g=1}^G N_g(N_g-1)} \sum_{g=1}^G \lambda_{1g}^* (\chi_1^2)_g + N \left(\mathbf{0}, \frac{16N_0^3}{\mu_2^2 \left[\sum_{g=1}^G N_g(N_g-1) \right]^2} (\gamma_1^*)' \Sigma_1^* \gamma_1^* \right) + \frac{2(\gamma_1^*)^2}{\mu_2}$$

where λ_{1g}^* are the characteristic roots of Σ_1^* .

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 $(\gamma_1^*)' \Sigma_1 = \mathbf{0}$ (Searle, 1971).

Recall that

$$\Sigma_1 = \begin{pmatrix} N_1(N_1-1)\tau_{11}^2 & \sqrt{N_1(N_1-1)N_2(N_2-1)}\tau_{112} & \dots & \sqrt{N_1(N_1-1)N_G(N_G-1)}\tau_{11G} \\ \sqrt{N_1(N_1-1)N_2(N_2-1)}\tau_{112} & N_2(N_2-1)\tau_{12}^2 & \dots & \sqrt{N_2(N_2-1)N_G(N_G-1)}\tau_{12G} \\ \vdots & \vdots & \ddots & \vdots \\ \sqrt{N_1(N_1-1)N_G(N_G-1)}\tau_{11G} & \sqrt{N_2(N_2-1)N_G(N_G-1)}\tau_{12G} & \dots & N_G(N_G-1)\tau_{1G}^2 \end{pmatrix}$$

where

$$\begin{aligned} \tau_{1g}^2 &= 4\xi_1^{(12)} \left\{ \frac{1}{N_g} + \frac{\sum_{g=1}^G N_g(N_g-1)^2}{M^2(M-1)^2} + \sum_{g < g'} \frac{N_g N_{g'}(N_{g'} + N_g)}{M^2(M-1)^2} + 2 \sum_{g \neq g' \neq g^\dagger} \frac{N_g N_{g'} N_{g^\dagger}}{M^2(M-1)^2} \right. \\ &\quad \left. + 2 \sum_{g=1}^G \sum_{g'=g+1}^G \frac{N_g(N_g-1)N_{g'}}{M^2(M-1)^2} - 2 \frac{(N_g-1)}{M(M-1)} - 2 \frac{(M-N_g)}{M(M-1)} \right\} \end{aligned}$$

and

$$\begin{aligned} \tau_{1gg'} &= \left(\frac{4\xi_1^{(12)}}{M} \right) \left[-2 + \frac{1}{M(M-1)^2} \left(\sum_{g=1}^G N_g(N_g-1)^2 + \sum_{1 \leq g < g' \leq G} N_g N_{g'}(N_g + N_{g'}) \right) \right. \\ &\quad \left. + 2 \sum_{\substack{1 \leq g, g', g^\dagger \leq G \\ g \neq g' \neq g^\dagger}} N_g N_{g'} N_{g^\dagger} + 2 \sum_{g=1}^G \sum_{g'=g+1}^G N_g(N_g-1)N_{g'} \right] \end{aligned}$$

Then, the first element of $(\gamma_1^*)' \Sigma_1$ is

$$\frac{\gamma_1^*}{\sqrt{N_0}} \sqrt{N_1(N_1-1)} [N_1(N_1-1)\tau_{11}^2 + \sum_{g=2}^G N_g(N_g-1)\tau_{11g}]$$

and

$$\begin{aligned} \frac{\gamma_1^*}{\sqrt{N_0}} \sqrt{N_1(N_1-1)} [N_1(N_1-1)\tau_{11}^2 + \sum_{g=2}^G N_g(N_g-1)\tau_{11g}] &= 0 \\ \Leftrightarrow N_1 = 1 \text{ or all } N_g \text{'s} = 1 \text{ or } (\tau_{11}^2 = 0 \text{ and all } \tau_{11g}) &= 0 \end{aligned}$$

So, $\mathbf{V}'_N \mathbf{V}_N$ and $(\gamma_1^*)' \mathbf{V}_N$ are independent if and only if $N_1 = 1$ or all N_g 's = 1 or $(\tau_{11}^2 = 0$ and all τ_{11g} 's) = 0, which is not the case here.

Now, write

$$[\mathbf{V}'_N \mathbf{V}_N + 2(\boldsymbol{\gamma}_1^*)' \mathbf{V}_N + (\boldsymbol{\gamma}_1^*)' \boldsymbol{\gamma}_1^*] = [(\mathbf{V}_N + \boldsymbol{\gamma}_1^*)' (\mathbf{V}_N + \boldsymbol{\gamma}_1^*)]$$

and

$$\begin{aligned} T_{N,2} &= \frac{2N_0}{\mu_2} (BMS) + O_p(N_0^{-1/2}) = \frac{2N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \mathbf{D}'_1 \mathbf{D}_1 + O_p(N_0^{-1/2}) \\ &= \frac{2N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} (\mathbf{V}_N + \boldsymbol{\gamma}_1^*)' (\mathbf{V}_N + \boldsymbol{\gamma}_1^*) + O_p(N_0^{-1/2}) \end{aligned}$$

Note that $(\mathbf{V}_N + \boldsymbol{\gamma}_1^*) \sim N(\boldsymbol{\gamma}_1^*, \boldsymbol{\Sigma}_1)$, $\left(\frac{\mathbf{V}_N}{\sqrt{N_0}} + \boldsymbol{\gamma}_1^*\right) \sim N(\boldsymbol{\gamma}_1^*, \boldsymbol{\Sigma}_1^*)$ and

$$\mathbf{D}_1 \sim N(\boldsymbol{\nu}_1, \boldsymbol{\Sigma}_1) \quad \text{or} \quad \frac{\mathbf{D}_1}{\sqrt{N_0}} \sim N(\boldsymbol{\gamma}_1^*, \boldsymbol{\Sigma}_1^*).$$

The distribution of $\frac{\mathbf{D}'_1 \mathbf{D}_1}{N_0}$ 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

$$\mathbf{Y} = \mathbf{P} \frac{\mathbf{D}_1}{\sqrt{N_0}} \Rightarrow \frac{\mathbf{D}_1}{\sqrt{N_0}} = \mathbf{P}'\mathbf{Y}$$

Then,

$$\mathbf{Y} \sim N(\mathbf{P}\boldsymbol{\gamma}_1^*, \boldsymbol{\Lambda}) \quad \text{and} \quad \frac{\mathbf{D}'_1 \mathbf{D}_1}{N_0} = \mathbf{Y}'\mathbf{P}\mathbf{P}'\mathbf{Y} = \mathbf{Y}'\mathbf{Y},$$

Hence,

$$\frac{\mathbf{D}'_1 \mathbf{D}_1}{N_0} = \mathbf{Y}'\mathbf{Y} \sim \sum_{i=1}^G \lambda_i (\chi_1^2(\delta_i)) \tag{7.1}$$

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 (7.1),

$$T_{N,2} = \frac{2N_0}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \mathbf{D}'_1 \mathbf{D}_1 \sim \frac{2N_0^2}{\mu_2 \sum_{g=1}^G N_g(N_g - 1)} \sum_{i=1}^G \lambda_i (\chi_1^2(\delta_i))_i$$

Since we have a linear combination of non-central chi-square random variables, when $\nu_1 = \frac{\gamma_1^*}{\sqrt{N_0}}$,

$$P(T_{N,2} > \nu_1) \rightarrow 1 \quad \text{as} \quad N_0 \rightarrow \infty$$

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 $T_{N,3}$.

8. Data analysis

The data set consists of three groups of HIV infected individuals. The interest is to compare the *env* gene V3 loops from B clade macrophage-tropic, B clade t-cell adapted and clade C sequences. There is a hypothesis which says that clade C is like B clade macrophage-tropic sequences. For this data set there are 356 independent sequences from B clade macrophage-tropic, 140 B clade t-cell adapted sequences and 151 clade C sequences. The sequences are all at the amino acid level with 35 positions long and they can be downloaded from the Los Alamos repository at the address <http://hiv-web.lanl.gov>

Since the elements of Σ_1 and Σ_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_{N_2} and T_{N_3} from the data set.
2. Sample 356 sequences for the B clade macrophage-tropic group, 140 sequences for B clade t-cell adapted group and 151 for clade C group with replacement from the pooled sample, i.e., the combined groups.
3. Recompute the test statistics T_{N_2} and T_{N_3} 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 $\frac{\#T'_{N_2}s \geq T_{N_2}obs}{R}$ and $\frac{\#T'_{N_3}s \geq T_{N_3}obs}{R}$.

The results are

$$T_{N_2}obs = 5.235 \quad T_{N_3}obs = 0.4173$$

For $R = 10,000$, the percentiles of the bootstrap distribution are given in Table 1.

Table 1: Percentiles of the Bootstrap Distribution

Statistic	1%	5%	95%	99%
T_{N_2}	1.275	1.290	1.544	1.574
T_{N_3}	0.399	0.404	0.457	0.459

When comparing the three groups, the observed p-value for T_{N_2} is less than 1/10001 and for T_{N_3} is 0.7586. Therefore, we can say that relative to the within-clade variation, there is significant variability between the three clades, but relative to the across-within-clade, there is no significant variability across-between the three clades.

Since we found a significant variability among the three clades, we now need to compare the variability of the groups two by two. Using the same procedure described above, we get:

For the comparison of B clade macrophage-tropic group and B clade t-cell adapted group,

$$T_{N_2}obs = 5.7975 \quad T_{N_3}obs = 2.9928$$

For $R = 20,000$, the percentiles of the bootstrap distribution are given in Table 2.

Table 2: Percentiles of the Bootstrap Distribution

Statistic	1%	5%	95%	99%
T_{N2}	1.5524	1.5954	1.7682	1.8006
T_{N3}	0.7929	0.8028	0.8463	0.8546

When comparing these two groups, the observed p-value for T_{N2} and T_{N3} are less than 1/20001. Therefore, we can say that relative to the within-clade variation, there is significant variability between the B clade macrophage-tropic group and the B clade t-cell adapted group and relative to the across-within-clade, there is significant variability across-between the two clades.

For the comparison of B clade macrophage-tropic group and clade C group,

$$T_{N2obs} = 4.6724 \quad T_{N3obs} = 2.3057$$

For $R = 20,000$, the percentiles of the bootstrap distribution are given in Table 3.

Table 3: Percentiles of the Bootstrap Distribution

Statistic	1%	5%	95%	99%
T_{N2}	1.4181	1.4331	1.4989	1.5130
T_{N3}	0.8964	0.9075	0.9579	0.9682

When comparing these two groups, the observed p-value for T_{N2} and for T_{N3} are less than 1/20001. Therefore, we can say that relative to the within-clade variation, there is significant variability between the B clade macrophage-tropic group and the clade C group. Also, relative to the across-within-clade, there is significant variability across-between the two clades.

Comparing B clade t-cell adapted group and clade C group,

$$T_{N2obs} = 5.0621 \quad T_{N3obs} = 1.1312$$

For $R = 20,000$, the percentiles of the bootstrap distribution are given in Table 4.

Table 4: Percentiles of the Bootstrap Distribution

Statistic	1%	5%	95%	99%
T_{N2}	3.2125	3.4611	4.9459	5.3518
T_{N3}	2.0066	2.1183	2.6653	2.7810

When comparing these two groups, the observed p-value for T_{N2} is 0.0335 and for T_{N3} is 1. Therefore, we can say that relative to the within-clade variation, there is significant variability between the B clade t-cell adapted group and the clade C group, but relative to the across-within-clade, there is no significant variability across-between the two clades.

References

- [Anderson and Landis, 1980] Anderson, R. J. and Landis, J. R. (1980). CATANOVA for multidimensional contingency tables: Nominal-scale response. *Communications in Statistics - Theory and Methods*, A9(11):1191–1206.
- [Anderson and Landis, 1982] Anderson, R. J. and Landis, J. R. (1982). CATANOVA for multidimensional contingency tables: Ordinal-scale response. *Communications in Statistics - Theory and Methods*, 11(3):257–270.
- [Hoeffding, 1948] Hoeffding, W. (1948). A class of statistics with asymptotically normal distribution. *Ann. of Math. Stat.*, 19:293–325.
- [Lee, 1990] Lee, A. J. (1990). *U-Statistics - Theory and Practice*. Marcel Dekker, Inc.
- [Light and Margolin, 1971] Light, R. J. and Margolin, B. H. (1971). An analysis of variance for categorical data. *Journal of the American Statistical Association*, 66:534–544.
- [Light and Margolin, 1974] Light, R. J. and Margolin, B. H. (1974). 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.
- [Mises, 1947] Mises, R. V. (1947). On the asymptotic distribution of the differentiable statistical functions. *Annals of Mathematical Statistics*, 18:309–348.
- [Pinheiro, 1997] Pinheiro, H. P. (1997). *Modelling Variability in the HIV Genome*. PhD thesis, University of North Carolina. Mimeo Series No. 2186T.
- [Pinheiro et al., 1999] Pinheiro, H. P., Seillier-Moiseiwitsch, F., and Sen, P. K. (1999). Multivariate CATANOVA and applications to DNA sequences in categorical data. Research Report 12/99, Universidade Estadual de Campinas.
- [Puri and Sen, 1971] Puri, M. L. and Sen, P. K. (1971). *Nonparametric Methods in Multivariate Analysis*. Wiley, New York.
- [Searle, 1971] Searle, S. R. (1971). *Linear Models*. John Wiley & Sons.
- [Searle, 1982] Searle, S. R. (1982). *Matrix Algebra Useful for Statistics*. John Wiley & Sons.
- [Seillier-Moiseiwitsch et al., 1994] Seillier-Moiseiwitsch, F., Margolin, B. H., and Swanstrom, R. (1994). Genetic variability of human immunodeficiency virus: Statistical and biological issues. *Annual Review of Genetics*, 28:559–596.
- [Sen, 1967] Sen, P. K. (1967). On some multisample permutation tests based on a class of u-statistics. *American Statistical Association Journal*, pages 1201–1213.
- [Sen, 1981] Sen, P. K. (1981). *Invariance Principles and Statistical Inference*. John Wiley & Sons.
- [Sen, 1995] Sen, P. K. (1995). Paired comparisons for multiple characteristics: An anocova approach. In Nagaraja, H. N., Sen, P. K., and Morrison, D. F., editors, *Statistical Theory and Practice: Papers in Honor of H. A. David*, pages 247–264. Springer-Verlag, New York.

[Sen and Singer, 1993] Sen, P. K. and Singer, J. M. (1993). *Large Sample Methods in Statistics*. Chapman & Hall.

[Weir, 1990] Weir, B. S. (1990). *Genetic Data Analysis*. Sinauer Associates.