Identities Emerging from Two Sample U-Statistics

SHARADA V. BHAT^{*} and I.D. SHETTY Department of Statistics Karnatak University, Dharwad

ABSTRACT

In this paper, we derive some mathematical identities which involve combinatorial coefficients. The well known theory of two-sample U Statistics is used to derive the identities.

Key words and Phrases: U-Statistics, two-sample location problem, combinatorial coefficients, ordered ranks.

1 Introduction

Identities are useful in simplifying many algebraic expressions. They provide simple alternate expressions to solve complex algebraic expressions. Riordan (1968) contains many such fundamental identities, Joshi and Balakrishnan (1981) provide statistical derivations of some such identities and Baiju and Thomas (2007) describe some identities using well established theories of order statistics and U-statistics based on certain linear functions of order statistics.

In this paper, we derive some identities using two-sample U-statistics and ordered ranks. The two-sample U-statistics is described in section 2 and two-sample

^{*}email: bhat_sharada@yahoo.com

U-Statistics from which these identities emerge are described in section 3. The identities derived from mean and variance of these U-statistics is given in section 4 and identities emerging from ordered ranks is given is section 5.

2 Definition of Two-sample U-Statistics

U-Statistics is a class of unbiased estimators of parameters of a population. They are often used as test statistics although they emphasize estimation. Randles and Wolfe (1979) describes the two-sample U-Statistics as follows:

Let X_1, \ldots, X_m and Y_1, \ldots, Y_n be the observations of two independent random samples drawn from cumulative distribution functions (cdf) F(x) and G(y) respectively. A parameter θ is said to be estimable of degree (b, d) for distributions (F, G)in a family \mathbf{F} if b and d are the smallest sample sizes for which there exists an estimator of θ that is unbiased for every $(F, G) \in \mathbf{F}$. That is, there is a function h(.; .)such that $E_{F,G}[h(X_1, \ldots, X_b; Y_1, \ldots, Y_d)] = \theta$ for every $(F, G) \in \mathbf{F}$, where h(.; .) is called as two-sample kernel and is symmetric in it's X_i components and separately symmetric in it's Y_j components. Under these conditions a two sample U-statistic, for $m \ge b$ and $n \ge d$ has the form

$$U(X_1, \dots, X_m; Y_1, \dots, Y_n) = \left(\begin{pmatrix} m \\ b \end{pmatrix} \begin{pmatrix} n \\ d \end{pmatrix} \right)^{-1} \sum_{\alpha} h(X_{i_1}, \dots, X_{i_b}, Y_{j_1}, \dots, Y_{j_d}) ,$$
(2.1)

where \sum_{α} is the collection of all subsets of b(d) integers chosen without replacement from the integers $\{1, \ldots, m\}$ $\{1, \ldots, n\}$.

3 Some two sample U-Statistics for location problem

Suppose X_1, \ldots, X_m and Y_1, \ldots, Y_n are independent random samples from populations with absolutely continuous distribution functions F(x) and G(y) having probability density functions (pdf) f(x) and g(y) respectively. Then the two sample location problem is to test $H_0: F(x) \equiv G(x)$ against the alternative $H_1: G(x) =$

 $F(x-\theta), \theta > 0$ or $\theta < 0$ or $\theta \neq 0, -\infty < x < \infty$, that is, two distributions differ only in their location parameter.

Suppose b and d are some fixed positive integers such that $1 \leq b \leq m$ and $1 \leq d \leq n$. For testing $H_0: \theta = 0$ against $H_1: \theta > 0$, Shetty and Bhat (1993) proposed

$$SB_{1} = \left[\binom{m}{b}\binom{n}{d}\right]^{-1} \sum_{\alpha} \phi_{1}\left(X_{i_{1}}, ..., X_{i_{b}}; Y_{j_{1}}, ..., Y_{j_{d}}\right),$$
(3.1)

where \sum_{α} is sum over all $\binom{m}{b}\binom{n}{d}$ possible sub samples,

$$\phi_1(X_1, ..., X_b; Y_1, ..., Y_d) = \begin{cases} 1, & if \quad M_1 \le M_2 \\ 0, & otherwise \end{cases} ,$$
 (3.2)

 M_1 = median of $(X_1, ..., X_b)$, M_2 = median of $(Y_1, ..., Y_d)$, and b and d are odd positive integers.

Shetty and Bhat (1994) proposed

$$SB_2 = [\binom{m}{b}\binom{n}{d}]^{-1} \sum_{\alpha} \phi_2(X_{i_1}, ..., X_{i_b}; Y_{j_1}, ..., Y_{j_d})$$
(3.3)

and

$$SB_{3} = \left[\binom{m}{b}\binom{n}{d}\right]^{-1} \sum_{\beta} \phi_{3}\left(X_{i_{1}}, ..., X_{i_{d}}; Y_{j_{1}}, ..., Y_{j_{b}}\right),$$
(3.4)

where \sum_{β} is sum over all $\binom{m}{d}\binom{n}{b}$ possible sub samples,

$$\phi_2(X_1, ..., X_b; Y_1, ..., Y_d) = \begin{cases} 1, & if \quad X_{(b)} \le M_2 \\ 0, & otherwise \end{cases}$$
(3.5)

$$\phi_3(X_1, ..., X_b; Y_1, ..., Y_d) = \begin{cases} 1, & if \quad M_3 \le Y_{(1)} \\ 0, & otherwise \end{cases}$$
(3.6)

 M_3 = median of $(X_1, ..., X_d)$, $X_{(b)}$ = maximum of $(X_1, ..., X_b)$, $Y_{(1)}$ = minimum of $(Y_1, ..., Y_b)$ and d is an odd positive integer.

Shetty et al.(1997) proposed a two-sample U-Statistic with kernel being the function of sample quantiles which is given by

$$SB_4 = [\binom{m}{b}\binom{n}{d}]^{-1} \sum_{\alpha} \phi_4(X_{i_1}, ..., X_{i_b}; Y_{j_1}, ..., Y_{j_d})$$
(3.7)

and

$$\phi_4(X_1, ..., X_b; Y_1, ..., Y_d) = \begin{cases} 1 & if \quad X_{(k_1)b} \le Y_{(k_2)d} \\ 0 & otherwise \end{cases}$$
(3.8)

where,

$$\mathbf{k}_{1} = \begin{cases} b\beta, & if \quad b\beta \quad \text{is an integer} \\ [b\beta] + 1, & if \quad b\beta \quad \text{is not an integer} \\ k_{2} = \begin{cases} d\gamma, & if \quad d\gamma \quad \text{is an integer} \\ [d\gamma] + 1, & if \quad d\gamma \quad \text{is not an integer} \\ \end{cases}$$

 $\beta^{\rm th}$ quantile of a sample size n is defined as the $r^{th} {\rm order}$ statistic, where

$$r = \begin{cases} n\beta, & if & n\beta & \text{is an integer} \\ [n\beta] + 1, & if & n\beta & \text{is not an integer} \end{cases}$$

 $X_{(k_1)b} = k_1^{th}$ order statistic of $(X_1, ..., X_b)$ and $Y_{(k_2)d} = k_2^{th}$ order statistic of $(Y_1, ..., Y_d)$. All these two sample U-Statistics are expressed in alternative forms using ordered ranks and their properties are also studied. An extensive study of these statistics is carried out in Bhat (1996).

4 Identities from mean and variance of two sample U-Statistics

In this section, we present some identities and their proofs based on mean and variance of two sample U statistics defined in section 3.

Identity 4.1.

$$\left(\frac{d!}{(q!)^2}\right) \sum_{i=p+1}^{b} {b \choose i} B(i+q+1,b-i+q+1) = 1/2$$
 (4.1)

for b and d being odd positive integers. **Proof.** The mean of SB_1 under H_0 is obviously 1/2. But it is also given by

$$\mathbf{E}(\mathbf{SB}_{1}) = \int_{-\infty}^{\infty} \mathbf{F}_{M_{1}}(x) \, \mathrm{d} \, \mathbf{F}_{M_{2}}(x) \, \mathrm{d} \, \mathbf{F}_{M_{2}}(x)$$

where

$$F_{M_1}(x) = \sum_{i=p+1}^{b} \begin{pmatrix} b \\ i \end{pmatrix} [F(x)]^i [\overline{F}(x)]^{b-i},$$

$$F_{M_2}(x) = \sum_{i=q+1}^{d} \begin{pmatrix} d \\ i \end{pmatrix} [F(x)]^i [\overline{F}(x)]^{d-i},$$

$$= (b-1)/2, \qquad q = (d-1)/2 \qquad \text{and} \qquad \overline{F}(x) = 1 - F(x) .$$

Therefore,

p

$$E(SB_{1}) = \int_{-\infty}^{\infty} \sum_{i=p+1}^{b} {\binom{b}{i}} [F(x)]^{i} [\overline{F}(x)]^{b-i}(x) (d!/(q!)^{2}) [F(x)]^{q} [\overline{F}(x)]^{q} dF(x)$$
$$= \sum_{i=p+1}^{b} {\binom{b}{i}} (d!/(q!)^{2}) \int_{-\infty}^{\infty} [F(x)]^{i+q} [\overline{F}(x)]^{b-i+q} dF(x)$$
$$= (d!/(q!)^{2}) \sum_{i=p+1}^{b} {\binom{b}{i}} B(i+q+1,b-i+q+1),$$

where $B(x, y) = \Gamma(x) \Gamma(y) / \Gamma(x + y)$. Therefore, we get the identity (4.1). Identity 4.2.

$$b^2 \zeta_{10}(SB_1) / \left(d^2 \zeta_{01}(SB_1) \right) = 1,$$
 (4.2)

where

$$\zeta_{10}(SB_1) = cov \left[\phi_1(X_1, ..., X_b; Y_1, ..., Y_d), \phi_1(X_1, X_{b+1}, ..., X_{2b-1}; Y_1, ..., Y_{2d})\right]$$

and

$$\zeta_{01}(SB_1) = cov \left[\phi_1(X_1,...,X_b;Y_1,...,Y_d) \ , \ \phi_1(X_{b+1},...,X_{2b};Y_1,Y_{d+1},...,Y_{2d-1})\right].$$

Proof. While deriving the asymptotic variance of SB_1 under the null hypothesis, certain expressions are evaluated. We have

$$\zeta_{10}(SB_1) = E\left[\phi_1(X_1, ..., X_b; Y_1, ..., Y_d), \phi_1(X_1, X_{b+1}, ..., X_{2b-1}; Y_{d+1}, ..., Y_{2d})\right] - E\left[SB_1\right]^2$$

$$= \int_{-\infty}^{\infty} P^2 \left(med \left(x, X_2, ..., X_b \right) \le M_2 \right) dF(x) - (1/4) .$$
(4.3)

Also

$$P(med (x, X_2, ..., X_b) \le M_2) = d!(b-i)![K_1 + \dots + K_b],$$
(4.4)

where

$$K_1 = P[x \le Y_1, X_2 \le X_3 \le \dots \le x \le \dots \le X_b, Y_2 \le Y_3 \le \dots \le Y_1 \le \dots \le Y_d],$$

 \boldsymbol{x} in the middle position and Y_1 in the middle position.

$$K_i = P[X_i \le Y_1, X_2 \le X_3 \le \dots \le X_i \le \dots \le x, Y_2 \le Y_3 \le \dots \le Y_1 \le \dots \le Y_d,]$$

for i = 2, ..., b and X_i in the middle position and Y_1 in the middle position.

After evaluating the expressions for K_1, \ldots, K_b and substituting in (4.4) and (4.3) we get

$$\zeta_{10}(SB_1) = (d!/(b-1)!)^2 K(b,d) / \left(p!^2 q!^2 (b+d-1) \begin{pmatrix} b+d-2\\ p+q \end{pmatrix} \right)^2, \quad (4.5)$$

where

$$K(b,d) = \sum_{i=0}^{p+q} {\binom{b+d-1}{i}}^2 B(2i+1, \ 2b+2d-2i-1) + \sum_{i\neq i'=1}^{p+q} \sum_{i=1}^{p+q} {\binom{b+d-1}{i}} {\binom{b+d-1}{i'}} B(i+i'+1, \ 2b+2d-i-i'-1) - (1/4)$$

Similarly,

$$\zeta_{01}(SB_1) = E[\phi_1(X_1, ..., X_b; Y_1, ..., Y_d) , \phi_1(X_{b+1}, ..., X_{2b} ; Y_1, Y_{d+1}, ..., Y_{2d-1}] - (1/4)$$

Identities by U-statistics

$$= (b!/(d-1)!)^{2} K(b,d) / \left(p!^{2}q!^{2}(b+d-1) \begin{pmatrix} b+d-2\\ p+q \end{pmatrix} \right)^{2}.$$
(4.6)

Therefore

$$\zeta_{10}(SB_1)/\zeta_{01}(SB_1) = (d!/(b-1)!)^2 / (b!/(d-1)!)^2 = (d/b)^2.$$
(4.7)

Hence, we get the identity (4.2).

Identity 4.3.

$$\zeta_{10} (SB_2) - \zeta_{01} (SB_3) = 0 , \qquad (4.8)$$

where

$$\zeta_{10}(SB_2) = cov \left[\phi_2(X_1, ..., X_b; Y_1, ..., Y_d), \phi_2(X_1, X_{b+1}, ..., X_{2b-1}; Y_{d+1}, ..., Y_{2d})\right]$$

and

$$\zeta_{01}(SB_3) = cov \left[\phi_3(X_1, ..., X_b; Y_1, ..., Y_d), \phi_3(X_{b+1}, ..., X_{2b}; Y_1, Y_{d+1}, ..., Y_{2d-1})\right].$$

Proof. Under the null hypothesis

$$E(SB_2) = \int_{-\infty}^{\infty} [F(x)]^b dF_{M_2}(x)$$

= $\left(\frac{d!}{(q!)^2}\right) B(b+q+1, q+1)$ (4.9)

and

$$E(SB_3) = \int_{-\infty}^{\infty} \left[\overline{F}(x)\right]^b dF_{M_3}(x),$$

where

$$F_{M_3}(x) = \sum_{i=q+1}^d \begin{pmatrix} d \\ i \end{pmatrix} [F(x)]^i \left[\overline{F}(x)\right]^{d-i}$$

Therefore,

$$E(SB_3) = \left(\frac{d!}{(q!)^2}\right) \int_{-\infty}^{\infty} \left[F(x)\right]^q \left[\overline{F}(x)\right]^{b+q} dF(x)$$

$$= (d!/(q!)^2) B (q+1, b+q+1) .$$
(4.10)

Since the kernel $\phi_2(.;.)$ can be obtained from $\phi_3(.;.)$ by replacing (X_i) 's by $(-X_i)$'s, (Y_j) 's by $(-Y_j)$'s and interchanging labels we get $E(SB_2) = E(SB_3)$.

Also,

$$\zeta_{10}(SB_2) = E\left(\phi_2(X_1, ..., X_b; Y_1, ..., Y_d), \phi_2(X_1, X_{b+1}, ..., X_{2b-1}; Y_{d+1}, ..., Y_{2d})\right) - E^2[SB_2]$$
$$= \int_{-\infty}^{\infty} P^2\left(\max\left(x, X_2, ..., X_b\right) \le M_2\right) dF(x) - E^2(SB_2)$$
(4.11)

and

$$\zeta_{01}(SB_3) = E\left(\phi_3(X_1, ..., X_d; Y_1, ..., Y_b), \phi_3(X_{d+1}, ..., X_{2d}; Y_1, Y_{b+1}, ..., Y_{2b-1})\right) - E^2[SB_3]$$
$$= \int_{-\infty}^{\infty} P^2\left(M_3 \le \min\left(y, Y_2, ..., Y_b\right)\right) dF(x) - E^2(SB_3). \tag{4.12}$$

Since $E[SB_2] = E[SB_3]$ under H_0 and from symmetry, we have

$$P(\max(x, X_2, ..., X_b) \le M_2) = P(M_3 \le \min(y, Y_2, ..., Y_d)).$$

From (4.9) through (4.12) we get the identity (??).

Identity 4.4.

$$k \binom{d}{k} \sum_{i=k}^{b} \binom{b}{i} B (i+k, 2b-i-k+1) = 1/2$$
(4.13)

 or

$$[d!/((k-1)!(d-k)!)] \sum_{i=k}^{b} {b \choose i} B (i+k, 2b-i-k+1) = 1/2.$$

Proof. Under the null hypothesis

$$E(SB_4) = P(X_{(k_i)} \le Y_{(k_2)})$$

= $\int_{-\infty}^{\infty} P(X_{(k_i)} \le y) dF_{Y_{(k_2)}}(y)$
= $k_2 \binom{d}{k_2} \sum_{i=k_1}^{b} \binom{b}{i} B(i+k_2, b+d-i-k_2+1)),$ (4.14)

When b = d and $k_1 = k_2 = k$, we have $E(SB_4) = 1/2$. Therefore we get identity (4.13).

Identity 4.5. For i = 1, 3, 5, ..., b + d - 1,

$$e(SB_1(b,d))/e(SB_1(i,b+d-i)) = 1,$$
(4.15)

where

$$e(SB_{1}(b,d)) = \left[(b!d!) / \left((p!q!)^{2} \sigma_{b,d} \right) \right] \int_{-\infty}^{\infty} [F(x)]^{p+q} \left[\overline{F}(x) \right]^{p+q} [f(x)]^{2} dx,$$

Identities by U-statistics

$$\sigma_{b,d}^{2} = (d!)^{2} (b!)^{2} K (b,d) / \left[(p!)^{4} (q!)^{4} (b+d-1)^{2} \begin{pmatrix} b+d-2\\ p+q \end{pmatrix}^{2} \lambda (1-\lambda) \right]$$

and

$$0 < \lambda = \lim_{N \to \infty} (m/N) < 1, \qquad N = m + n$$

Proof. Under the null hypothesis

$$\sigma_{b,d}^{2} = \left(b^{2}/\lambda\right)\zeta_{10}\left(SB_{1}\right) + \left(d^{2}/\left(1-\lambda\right)\right)\zeta_{01}\left(SB_{1}\right)$$

$$= b^{2}\zeta_{10}\left(SB_{1}\right) / \left(\lambda\left(1-\lambda\right)\right) \quad \text{or} \quad d^{2}\zeta_{01}\left(SB_{1}\right) / \left(\lambda\left(1-\lambda\right)\right) \quad \text{by identity 4.2.}$$

It is worth to note that $e(SB_1(b,d))$ depends on (b+d) and underlying distribution F(x). Thus for $1 \le b \le m$, $1 \le d \le n$, b,d being odd positive integers, given F(x), we get $e(SB_1(b,d)) = e(SB_1(i,b+d-i))$ for i = 1, 3, 5, ..., b+d-1. Therefore, we get identity (4.15).

5 Identities based on Ordered Ranks

In this section, we present some identities based on the ordered ranks of two sample U statistics defined in section 3. Suppose that $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(m)}$ and $Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(n)}$ are the order statistics corresponding to X and Y sample observations respectively. Let $R_{(i)}(S_{(j)})$ be the rank of $X_{(i)}(Y_{(j)})$ in the joint ranking of X and Y observations. Then we have the following identities.

Identity 5.1.
$$\sum_{i=1}^{m} \sum_{j=0}^{q} \binom{i-1}{p} \binom{m-i}{p} \binom{R_{(i)}-i}{q-j} \binom{n-R_{(i)}+i}{q+1+j}$$
$$= \sum_{j=1}^{n} \sum_{i=0}^{p} \binom{j-1}{q} \binom{n-j}{q} \binom{S_{(j)}-j}{p+1+i} \binom{m-S_{(j)}+j}{p-i}.$$
(5.1)

Proof. Choose a sub sample of size *b* from the *X* sample such that $X_{(i)}$ is the median. For a fixed *i*, this can be done in $\binom{i-1}{p}\binom{m-i}{p}$ ways. Similarly choose a sub sample of size *d* from *Y* sample such that $Y_{(j)}$ is the median and is

greater than $X_{(i)}$. Each such sub-sample pair results in $\phi_1(.;.) = 1$. Using the fundamental rules of counting, we get the identity (5.1).

Identity 5.2.

$$\binom{m}{b}\binom{n}{d}SB_2 = \sum_{j=1}^n \binom{j-1}{q}\binom{n-j}{q}\binom{S_{(j)}-j}{b}.$$
 (5.2)

Proof. Choose a sub sample of size d from the Y sample such that $Y_{(j)}$ is the median. For a fixed j this can be done in $\binom{j-1}{q}\binom{n-j}{q}$ ways. The number of X observations less than $Y_{(j)}$ will be $(S_{(j)} - j)$. A sub sample of size b from the X observations can be chosen in $\binom{S_{(j)} - j}{b}$ ways and we get identity (5.2).

Identity 5.3.

$$\binom{m}{d}\binom{n}{b}SB_3 = \sum_{i=1}^m \binom{i-1}{q}\binom{m-i}{q}\binom{n-R_{(i)}+i}{b}.$$
 (5.3)

Proof. For a fixed*i*, $X_{(i)}$ can be chosen as median of sub sample of size *d* from Xobservations in $\binom{i-1}{q}\binom{m-i}{q}$ ways. The number of Yobservations greater than $X_{(i)}$ is $(n - R_{(i)} + i)$. A sub sample of size *b* from these *Y* observations can be chosen in $\binom{n-R_{(i)}+i}{b}$ ways. For each *i*, $\binom{i-1}{q}\binom{m-i}{q}\binom{n-R_{(i)}+i}{b}$ sub-sample pairs for which $\phi_3(.;.) = 1$. Then by the fundamental rule of counting, we get the identity (5.3).

Identity 5.4.

$$\sum_{i=1}^{m} \sum_{j=0}^{k_{2}-1} \binom{i-1}{k_{1}-1} \binom{m-i}{b-k_{1}} \binom{R_{(i)}-i}{k_{2}-j-1} \binom{n-R_{(i)}+i}{d-k_{2}+j+1} = \sum_{j=1}^{n} \sum_{i=0}^{k_{1}-1} \binom{j-1}{k_{2}-1} \binom{n-j}{d-k_{2}} \binom{S_{(j)}-j}{b-k_{1}+i+1} \binom{m-S_{(j)}+j}{k_{1}-i-1}.$$
(5.4)

Proof. Choose a sub sample of size *b* from the *X* sample such that $X_{(i)}$ is the k_1^{th} order statistic $((b\beta)^{th}$ quantile). For a fixed *i*, this can be done in $\begin{pmatrix} i-1\\ k_1-1 \end{pmatrix}$

 $\begin{pmatrix} m-i \\ b-k_1 \end{pmatrix}$ ways. Now choose a sub sample of size d from Y sample such that $Y_{(j)}$ is the k_2^{th} order statistic $((d\gamma)^{th}$ quantile) and is greater than $X_{(i)}$. Each such sub-sample pair results in $\phi_4(.;.) = 1$. The $Y_{(j)}$ can be selected in $\begin{pmatrix} R_{(i)} - i \\ k_2 - j - 1 \end{pmatrix} \begin{pmatrix} n - R_{(i)} + i \\ d - k_2 + j + 1 \end{pmatrix}$ ways. Thus using the fundamental rules of counting, we get the identity (5.4).

Acknowledgement: We thank the referee for his useful comments.

References

- Baiju, K.V. and Thomas, P.Y. (2007). On some identities proved using Statistical mathematics, Journal of the Kerala Statistical Association 18, 34-44.
- Bhat, S.V. (1996). Studies in Nonparametric Inference. An unpublished Ph.D thesis submitted to Karnatak University, Dharwad.
- Joshi, P.C. and Balakrishnan, N. (1981). Applications of order Statistics in combinatorial identities, Journal of Combinatorics, Information and System Sciences 6, 271-278.
- Randles, R. H. and Wolfe, D.A. (1979). Introduction to Theory of Nonparametric Statistics, John Wiley and Sons, New York.
- Riordan, J. (1968). Combinatorial Identities, John Wiley and Sons, New York.
- Shetty, I.D. and Bhat, S.V. (1993). Some Competitors of Mood's median test for location alternatives, Karnatak University, Journal of Science 37, 138-146.
- Shetty, I.D. and Bhat, S.V. (1994). A Note on the Generalization of Mathisen's median test, Statistics and Probability Letters 19, 199-204.
- Shetty, I.D, Sengupta, D. and Bhat, S. V (1997). A General class of two-sample Nonparametric tests based on sub-sample quantiles, Karnatak University, Journal of Science 41, 104-125.