# THE TREATMENT OF TIES IN THE WILCOXON TEST[1]

By Wolfgang J. Bühler

*University of California, Berkeley*

**1. Introduction.** Let $(X_1, \cdots, X_n)$ be a sample of $n$ independent observations from a distribution $F$, and $(Y_1, \cdots, Y_m)$ be a sample of independent observations from $G$. Then, if all $m + n$ observations are different, the Wilcoxon test will reject the hypothesis $F = G$, when the sum $S_{nm}$ of the ranks $R_i$ of the $X_i$ is too small or too large.

For the case with a positive probability of ties two procedures have been proposed. One is to order the tied observations randomly, the other is to replace $S_{nm}$ by $S'_{nm} = \sum_{i=1}^{n} R_i'$. Here $R_i' = \text{midrank}(X_i) = \frac{1}{2}[N_1(i) + N_2(i) + 1]$. $N_1(i)$ is the number of observations smaller than $X_i$ and $N_2(i)$ is the number of observations (including $X_i$) not larger than $X_i$.

If there are only finitely many values $\xi_k$ at which ties may occur and if $p_k = P\{X_1 = \xi_k\}$, then as shown by Putter [3] under certain regularity conditions the asymptotic relative efficiency of the "randomized" with respect to the midrank test is $1 - \sum_{k=1}^{n} p_k^3$. Using a slight modification of Putter's argument this note will show that this conclusion is still true if $p_k = P\{X_1 = \xi_k\} > 0$ and $q_k = P\{Y_1 = \xi_k\} > 0$ for infinitely many values $\xi_k$. The result is illustrated by applying it to certain parametric families of distributions, for which the efficiency of the midrank test has been investigated by Chanda [1]. Putter's notation will be used throughout the paper.

**2. The basic theorem.** Following Putter, let for

$$k = 1, 2, \cdots, \quad p_k = P\{X_1 = \xi_k\} > 0, \qquad q_k = P\{Y_1 = \xi_k\} > 0;$$

$U_k = $ number of $X$'s equal to $\xi_k$, $V_k = $ number of $Y$'s equal to $\xi_k$; $U = (U_1, U_2, \cdots), V = (V_1, V_2, \cdots), W = U + V; S_{nm}^0 = $ any statistic whose distribution is that of $S_{nm}$ under $F = G$; $\mu_{nm} = ES_{nm}^0 = n(n + m + 1)/2$, $\sigma_{nm}^2 = \text{Var } S_{nm}^0 = nm(n + m + 1)/12; T_{nm}^0 = (S_{nm}^0 - \mu_{nm})/\sigma_{nm}$.

Then the following theorem connects the asymptotic distributions of $S_{nm}$ and of $S'_{nm}$.

THEOREM 1. *If $m/n$ converges to a positive number $c$ as $m, n \to \infty$, then we have for any pair $(F, G$, of distributions with common discontinuities $\xi_k$, $k = 1, 2, \cdots$*

$$(2.1) \qquad \sigma_{U_k V_k}^2 / \sigma_{nm}^2 = a_k^2 \to_P b_k^2 = (1 + c)^{-1} p_k q_k(\theta)[p_k + c q_k(\theta)]$$

$$(2.2) \quad (S_{nm} - ES_{nm})/\sigma_{nm} = T_{nm} \to_{\mathcal{L}} N(0, b^2)$$

$$(2.3) \quad (S'_{nm} - ES_{nm})/\sigma_{nm} = T'_{nm} \to_{\mathcal{L}} N(0, \bar{b}^2),$$