

Revista Colombiana de Estadística
 Vols. 21-22, 1990

PRUEBAS INSEGADAS BASADAS EN RACHAS
 PARA ALTERNATIVAS DE LOCALIZACION Y ESCALA

Jimmy A. Corzo S.

Profesor Asistente
 Universidad Nacional

Resumen. En este artículo se presentan varias generalizaciones de la prueba de Wald Wolfowitz (1940) y de las pruebas para localización y escala propuestas por Ortiz (1983), Corzo & Ortiz (1983), Fernández & Ortiz (1986). Se establece una relación funcional entre rachas y rangos a través de la cual se demuestra que una subfamilia de las pruebas aquí propuestas es insesgada. Se presenta también una transformación de las estadísticas de prueba que incluyen sumas con número aleatorio de sumandos en sumas con número fijo de sumandos lo cual facilita la demostración de algunas de sus propiedades.

1. Preliminares.

La sucesión de variables aleatorias (v.a.'s) X_i , $i = 1, \dots, m, m+1, \dots, N$, $N = m+n$, $m, n \in \mathbb{N}$ independientes definidas sobre un espacio de probabilidad $(\Omega, \mathcal{A}, \mathcal{P})$ se llama una muestra aleatoria (m.a.); denótese por $X_{1;N}, \dots, X_{N;N}$ las corres

pondientes estadísticas de orden y sea F la Función de Distribución continua de X_i , $i = 1, \dots, N$.

Se denota por R_i el rango de X_i y por D_i el antirango de $X_{j;N}$, $i = 1, \dots, N$, y se definen por medio de las ecuaciones:

$$X_i = X_{R_i;N}, \quad i = 1, \dots, N$$

(1.1)

$$X_{D_j} = X_{j;N}, \quad j = 1, \dots, N.$$

Entonces R_i es el número de orden de X_i , $i = 1, \dots, N$ mientras que D_j es el subíndice de la v.a. de donde viene $X_{j;N}$, $j = 1, \dots, N$. En otras palabras: los eventos $(R_i = k)$ y $(D_j = t)$ indican respectivamente que X_i es la k -ésima v.a. más pequeña $i = 1, \dots, N$ y que $X_{j;N}$ es la estadística de orden correspondiente a la t -ésima v.a.

Por medio de la sucesión de v.a.

$$(1.2) \quad \eta_{R_i} = \begin{cases} 1, & i = 1, \dots, m \\ 0, & i = m+1, \dots, N, \quad i = 1, \dots, N \end{cases}$$

quedan representadas las primeras m v.a.'s de la muestra con unos y las últimas n v.a.'s de la muestra con ceros. La sucesión η_{R_i} , $i = 1, \dots, N$ se llama muestra dicotomizada.

Se definen las v.a.'s contadores

$$I_{R_i} = \begin{cases} 1, & \eta_{R_{i-1}} \neq \eta_{R_i} \\ 0, & \eta_{R_{i-1}} = \eta_{R_i}, \quad i = 1, \dots, N \end{cases}$$

con la condición de que $I_{R_i} = 1$ cuando $X_i = \min\{X_1, \dots, X_N\}$

es decir cuando $R_i = 1$.

Entonces de (1.2) y (1.3) se define el número de rachas hasta la v.a. X_i así

$$(1.4) \quad n_{R_i} = \sum_{\{k: R_k \leq R_i\}} I_{R_k} \quad i = 1, \dots, N.$$

En particular si $R_i = N$ es porque $X_i = \max\{X_1, \dots, X_N\}$ y por esto se puede interpretar a n_N como el número de rachas hasta la v.a. $\max\{X_1, \dots, X_N\}$. En general nótese que se puede interpretar como el número de rachas hasta la v.a. que se encuentra en la posición R_i en la sucesión X_1, \dots, X_N .

Entonces de la relación $\{R_i = j\} \Leftrightarrow \{D_j = i\}$ y puesto η_{R_i} , I_{R_i} y n_{R_i} son medibles respecto a $\sigma(R_i)$ (la σ -álgebra generada por R_i , $i = 1, \dots, N$), se obtiene para $j = 1, \dots, N$.

$$(1.5) \quad \eta_j = E(\eta_{R_i} / R_i = j) \begin{cases} 1, & \{1 \leq D_j \leq m\} = \bigcup_{i=1}^m \{R_i = j\} \\ 0, & \{m+1 \leq D_j \leq N\} = \bigcup_{i=m+1}^N \{R_i = j\} \end{cases}$$

$$(1.6) \quad I_j = E(I_{R_i} / R_i = j) = \begin{cases} 1, & \eta_{j-1} \neq \eta_j \\ 0, & \eta_{j-1} = \eta_j \end{cases}$$

con $I_i = 1$ y en consecuencia

$$\begin{aligned} n_1 &= E(n_{R_i} / R_i = j) = \sum_{k=1}^j I_k, \quad j = 1, \dots, N \\ &= 1 + \sum_{k=2}^j I_k, \quad j = 2, \dots, N. \end{aligned}$$

Nótese que r_{R_i} indica el número de rachas hasta la observación que se encuentra en la posición R_i , $i = 1, \dots, N$, mientras que r_j es el número de rachas hasta la j -ésima estadística de orden, $j = 1, \dots, N$. Pero por (1.1) $X_{j;N}$ es la j -ésima estadística de orden si $R_i = j$, es decir si X_i se encuentra en la posición $R_i = j$. Este argumento demuestra que las sucesiones r_{R_i} , $i = 1, \dots, N$ y r_j , $j = 1, \dots, N$ son equivalentes en el sentido que ambas contienen la misma información. Por tanto aquí se hace referencia a cualquiera de ellas con el nombre de muestra dicotomizada y se usarán convenientemente de acuerdo con lo que se quiere ilustrar.

Esta muestra dicotomizada tiene la siguiente estructura

$$(1.8) \quad \eta_1 = \dots = \eta_{\lambda_1} \neq \eta_{\lambda_1+1} = \dots = \eta_{\lambda_2} \neq \dots \neq \eta_{\lambda_{r_N-1}+1} = \dots = \eta_N$$

donde

$$(1.9) \quad \lambda_j = \sum_{i=1}^j L_i, \quad j = 1, \dots, r_N$$

en la cual hay $r_N \geq 2$ grupos de unos y ceros. Cada uno de estos grupos se denomina una racha y el número de elementos en la j -ésima racha L_j , $j = 1, \dots, r_N$ se llama longitud de la racha. r_N es el número total de rachas en la muestra dicotomizada. Obviamente r_j y L_j , $j = 1, \dots, r_N$ son v.a's.

Observación 1. λ_j , $j = 1, \dots, N$ representa la longitud acumulada de las primeras j rachas. Además para cada $j = 1, \dots, N$, λ es un tiempo de parada respecto a la sucesión de σ -álgebras $\mathcal{D}_{L_j} = \sigma(\mathcal{D}_1, \dots, \mathcal{D}_{L_j})$, $j = 1, \dots, r_N$. Corzo (1989).

De (1.8) y (1.9) se deducen también las siguientes relaciones

$$(1.10) \quad I_1 = \dots = I_{\lambda_1+1} = I_{\lambda_2+1} = \dots = I_{\lambda_{n_N-1}+1} = 1$$

$$(1.11) \quad I_{\lambda_{j-1}+2} = I_{\lambda_{j-1}+3} = \dots = I_{\lambda_j} = 0 \quad j=1, \dots, n_N$$

Por otra parte de (1.7) se puede calcular

$$\begin{aligned} n_{\lambda_{j-1}+k} &= \sum_{i=1}^{\lambda_{j-1}+k} I_i \\ &= \sum_{i=1}^{L_1} I_i = \sum_{i=1}^{L_2} I_{\lambda_1+i} + \dots + \sum_{i=1}^{L_{j-1}} I_{\lambda_{j-2}+i} + \sum_{i=1}^k I_{\lambda_{j-1}+i} \end{aligned}$$

entonces reemplazando las relaciones (1.10) y (1.11) en la anterior ecuación se consigue la igualdad

$$(1.12) \quad n_{\lambda_{j-1}+1} = \dots = n_{\lambda_j} = j$$

y de aquí sigue la relación

$$jL_j = \sum_{k=1}^{L_j} j = \sum_{k=1}^{L_j} n_{\lambda_{j-1}+k}.$$

Igualmente para cualquier función real g vale la igualdad.

$$g(j)L_j = \sum_{k=1}^{L_j} g(n_{\lambda_{j-1}+k})$$

Esta igualdad y el hecho de que

$$\sum_{i=1}^{n_N} L_i = N.$$

permiten establecer la siguiente transformación

$$(1.14) \quad \sum_{j=1}^{r_N} g(j) L_j = \sum_{j=1}^{r_N} \sum_{k=1}^{L_j} g(r_{\lambda_{j-1}} + k) \\ = \sum_{i=1}^N g(r_i).$$

Observación 2. La ecuación (1.14) se utilizará más adelante para transformar una familia de estadísticas de prueba que originalmente es una suma con número aleatorio de sumandos, en una suma con número fijo de sumandos.

2. Problema de dos Muestras.

Se considera aquí las v.a's. independientes X_1, \dots, X_m y X_{m+1}, \dots, X_N , con $N = m+n$ y $m, n = 1, 2, \dots$, sobre un espacio de probabilidad (Ω, \mathcal{A}, P) . Las primeras m de estas v.a's. se denominan muestra 1 y las restantes n constituyen la muestra 2. La muestra 1 tiene función de distribución continua F y la muestra 2 tiene función de distribución continua G . A la sucesión completa $X_1, \dots, X_m, X_{m+1}, \dots, X_N$, se le llama muestra combinada. También se considera aquí la muestra dicotomizada η_1, \dots, η_N que se obtiene como en (1.5).

En este párrafo se consideran las siguientes hipótesis:

Hipótesis nula

$$(2.1) \quad H_0 : F(x) = G(x), \quad \forall x \in \mathcal{R}$$

Alternativa de localización de una cola.

$$(2.2) \quad K_1 : F(x) = G(x-\mu), \quad \mu < 0, \quad \forall x \in \mathbb{R}$$

o bien $\mu > 0, \quad \forall x \in \mathbb{R}.$

Alternativa de localización de dos colas.

$$(2.3) \quad K_2 : F(x) = G(x-\mu), \quad \mu \neq 0, \quad \forall x \in \mathbb{R}.$$

Sean μ_1 y μ_2 las medianas de F y G respectivamente y $\theta = \mu_1 - \mu_2$.

Alternativa de escala de una cola.

$$(2.4) \quad K_3 : F(x) = G\left(\frac{x-\theta}{\zeta}\right), \quad \zeta \in (0,1), \quad \forall x \in \mathbb{R}$$

o bien $\zeta \in (1,\infty) \quad \forall x \in \mathbb{R}.$

Alternativa de escala de dos colas.

$$(2.5) \quad K_4 : F(x) = G\left(\frac{x-\theta}{\zeta}\right), \quad \zeta \in (0,\infty) \text{ y } \zeta \neq 1 \quad \forall x \in \mathbb{R}.$$

Se define la estadística de prueba para el problema de dos muestras como sigue:

$$(2.6) \quad T^* = \sum_{i=1}^{r_N} \delta_i^* h(i, L_i, r_N)$$

donde

$$\delta_i^* = \begin{cases} 1/m, & \text{si } \eta_{\lambda_{i-1}+1} = \dots = \eta_{\lambda_i} = 1 \\ -1/n, & \text{si } \eta_{\lambda_{i-1}+1} = \dots = \eta_{\lambda_i} = 0 \end{cases}$$

y h es una función real, tal que para algún $M \in \mathbb{R}$

$$0 < h(i, L_i, r_N) < M.$$

Además la función h se escoge creciente en i y L_i para pruebas con alternativas de localización y creciente en $\left| i - \frac{r_{N+1}}{2} \right|$ y en L_i , para pruebas con alternativas de escala.

Nótese que

$$(2.7) \quad |T^*| < r_N M < NM.$$

Observación 3. La distribución exacta de δ_i^* así como una manera de obtener la función h a partir de esperanzas condicionales se pueden consultar en Corzo (1989).

(i) Pruebas para Alternativas de Localización.

En este caso se escoge h creciente en i y L_i . Entonces como δ_i^* pondera los valores de h para la muestra 1 con $1/m$, los valores de T^* tienden a ser positivos cuando hay rachas grandes de la muestra 1 en las últimas posiciones de la muestra dicotomizada η_1, \dots, η_N . Recíprocamente cuando se presentan rachas grandes de la muestra 2 en las últimas posiciones de la muestra dicotomizada η_1, \dots, η_N , los valores de T^* tienden a ser negativos pues en este caso los valores de h se ponderan con $-1/n$. Es decir T^* toma valores positivos cuando $\mu > 0$ y valores negativos cuando $\mu < 0$ lo cual se puede interpretar como una prueba aleatorizada de nivel $\alpha \in (0, 1)$ para probar H_0 vs. K_1 con $\mu > 0$ así:

$$(2.8) \quad \phi_1 = \begin{cases} 1, & \text{si } T^* > c \\ \gamma, & \text{si } T^* = c \\ 0, & \text{si } T^* < c \end{cases}$$

donde para P_0 la distribución de T^* bajo H_0 los valores de c y γ se obtienen de las fórmulas siguientes:

$$c = \inf\{t : P_0(T^* > t) \leq \alpha\}$$

y

$$\gamma = \begin{cases} \frac{\alpha - P_0(T^* > c)}{P_0(T^* = c)}, & P_0(T^* = c) \neq 0 \\ 0 & , P_0(T^* = c) = 0 \end{cases}$$

Análogamente para la prueba de H_0 vs. K_2 :

$$(2.9) \quad \phi_2 = \begin{cases} 1, & \text{si } T^* < c_1, \quad T^* > c_2 \\ \gamma_i, & \text{si } T^* = c_i, \quad i = 1, 2 \\ 0, & \text{si } c_1 < T^* < c_2 \end{cases}$$

donde para P_0 la distribución de T^* bajo H_0 los valores de c y γ se obtienen de las fórmulas siguientes:

$$c_i = \inf\{t_i : P_i \leq \alpha/2\},$$

donde $P_1 = P_0(T^* < t_1)$ y $P_2 = P_0(T^* > t_2)$

y

$$\gamma_i = \begin{cases} \frac{\alpha/2 - P_i}{P_0(T^* > c_i)}, & P_0(T^* = c_i) \neq 0 \\ 0 & , P_0(T^* = c_i) = 0, \quad i = 1, 2 \end{cases}$$

(ii) Pruebas para Alternativas de Escala

En este caso se escoge h creciente en $\left| i - \frac{N+1}{2} \right|$ y L_i . Entonces como δ_i^* pondera los valores de h para la muestra 1 con $1/m$, los valores de T^* tienden a ser positivos cuando hay rachas grandes de la muestra 1 en los dos extremos de la muestra dicotomizada η_1, \dots, η_N . Recíprocamente cuando se presentan rachas grandes de la muestra 2 en los dos extremos de la muestra dicotomizada η_1, \dots, η_N los valores de T^* tienden a ser negativos pues en este caso los valores de h se ponderan con $-1/m$. Es decir T^* toma valores positivos cuando $\zeta \in (1, \infty)$ y valores negativos cuando $\zeta \in (0, 1)$ lo cual se puede interpretar como una prueba aleatorizada de nivel $\alpha \in (0, 1)$ para probar H_0 vs. K_3 con $\zeta \in (0, 1)$ análoga a (2.8) y de la misma manera se obtiene una prueba para H_0 vs. K_4 análoga a (2.9).

(iii) Casos especiales.

Los casos especiales se obtienen de acuerdo con la manera como se escoja la forma de la función h y las ponderaciones δ^* .

Por ejemplo para

$$(2.10) \quad h(i, L_i, r_N) = \begin{cases} m & 1 \leq D_i \leq m \\ -n & m+1 \leq D_i \leq N \end{cases}.$$

se obtiene $T^* = r_N$ que es la estadística propuesta por Wald y Wolfowitz (1940) para probar H_0 vs. la alternativa general; de dos colas $F(x) \neq G(x)$. Esta prueba se conoce con el nombre de wald y wolfowitz.

Observación 4. Wald y Wolfowitz mostraron en su artículo de 1940 que cuando $\frac{m}{n} \rightarrow c$, $c \neq 0$ y $c < \infty$, la prueba basada en \mathcal{L}_N es consistente. Más tarde Lehman (1951) sin demostrarlo comentaba que la prueba basada en \mathcal{L}_N podría ser insesgada. En este artículo (§5) se demuestra que no solo \mathcal{L}_N sino toda una subfamilia de pruebas basadas en T^* es insesgada. Posteriormente Mood (1954) demostró que la eficiencia asintótica relativa de la prueba de Wald y Wolfowitz para alternativas de localización y escala, en comparación con las conocidas pruebas t y F bajo el supuesto de muestreo de la distribución Normal, es igual a cero. Este resultado no es sorprendente si se tiene en cuenta que \mathcal{L}_N solo contiene el número de rachas hasta $X_{(N)} = \max\{X_1, \dots, X_N\}$ y por esto no contiene información sobre los agrupamientos de rachas en los extremos que son los que detectan las diferencias en localización y escala.

Otro caso especial de T^* es cuando $m = n$ en el cual es suficiente con la ponderación

$$\delta_i^* = \begin{cases} 1, & \text{si } \eta_{\lambda_{i-1}+1} = \dots = \eta_{\lambda_i} = 1 \\ -1, & \text{si } \eta_{\lambda_{i-1}+1} = \dots = \eta_{\lambda_i} = 0 \end{cases}$$

porque δ_i^* , $i = 1, \dots, \mathcal{L}_N$ en (2.6) tenía que ser dividido por m y n para neutralizar el efecto de diferencias grandes entre m y n pues podrían aparecer rachas largas en los extremos de η_1, \dots, η_N de aquella muestra cuyo tamaño es mayor sin que esto necesariamente indique diferencias en localización o escala. Además esta normalización es adecuada porque elimina el efecto de las diferencias en tamaños de muestra en el sentido que

$$\frac{1}{m} \sum_{i=1}^{r^{(1)}} L_i^{(1)} = \frac{1}{n} \sum_{i=1}^{r^{(0)}} L_i^{(0)} = 1$$

donde $L_i^{(1)}$ es la longitud de la i -ésima racha de unos, y $r^{(1)}$ es el número de rachas de unos en la muestra dicotomizada η_1, \dots, η_N (análogamente se definen $L_i^{(0)}$ y $r^{(0)}$).

Tomando δ_i^* como en (2.11) y los dos casos especiales de función h de (2.12) y (2.13) se obtienen las estadísticas propuestas por Ortiz (1983) y Corzo y Ortiz (1983) para dos muestras de igual tamaño ($m = n$), respectivamente:

$$(2.12) \quad h(i, L_i, r_N) = \frac{iL_i}{r_{N-1}}$$

$$(2.13) \quad h(i, L_i, r_N) = \frac{1}{r_{N-1}} \left[i - \frac{r_{N+1}}{2} \right]^2$$

mientras que son δ_i^* como en (2.6) y h de (2.12) se obtiene la propuesta de Fernández y Ortiz (1986). Además escogiendo δ_i^* como en (2.6) y h de (2.13) se obtiene la correspondiente estadística para el problema de dos muestras de diferentes tamaños con alternativa de escala.

El último caso especial que se presenta aquí es cuando la función h toma la forma

$$h(i, L_i, r_N) = Q(i, L_i)$$

es decir que no depende explícitamente de r_N que es la familia de estadísticas propuestas por Ortiz (1983). Al respecto cabe anotar que aunque la correlación entre las longitudes y el total de rachas r_N es negativa (a medida que aumen

tan las longitudes de las rachas disminuye su número), no parece ser muy alta pues las longitudes de una o todas las rachas pueden cambiar sin que se altere el valor de λ_N . Por esta razón no parece muy plausible eliminar de h el total de rachas en la muestra dicotomizada como propone Ortiz en su artículo.

3. Problema de una Muestra y Muestras Pareadas.

En este caso la muestra esta constituida por las v.a's. independientes X_1, \dots, X_N sobre un espacio de probabilidad $(\Omega, \mathcal{A}, \mathcal{P})$ con función de distribución continua F y con mediana

$$(3.1) \quad \mu = F^{-1}(\frac{1}{2}) .$$

Se consideran aquí las siguientes hipótesis:

Hipótesis nula

$$(3.2) \quad H'_0 : \mu = \mu_0$$

Alternativa de una cola y dos colas

$$(3.3) \quad K'_1 : \mu > \mu_0 \quad \text{o bien} \quad \mu < \mu_0$$

$$K'_2 : \mu \neq \mu_0$$

con $\mu_0 \in \mathbb{R}$ conocido.

La dicotomización de la muestra se obtiene a través de las v.a's. ξ_1, \dots, ξ_N , donde

$$\varepsilon_i = \begin{cases} 1 & \text{si } X_i > \mu_0 \\ 0 & \text{si } X_i < \mu_0 \end{cases} .$$

El número de observaciones mayores que μ_0 :

$$(3.5) \quad \varepsilon = \sum_{i=1}^N \varepsilon_i$$

es una v.a. que tiene distribución Binomial con parámetros N y $\theta = P(X_i > \mu_0)$. Bajo H_0 $\theta = \frac{1}{2}$ y por tanto

$$(3.6) \quad \lim_{N \rightarrow \infty} P(\varepsilon = 0) = \lim_{N \rightarrow \infty} P(\varepsilon = N) = \lim_{N \rightarrow \infty} \left(\frac{1}{2}\right)^N = 0$$

Por otra parte sea $\varepsilon_{\lambda_{j-1}^{(i)}+k}$ el k -ésimo elemento de J -ésima racha de ies, $i = 0, 1$ y λ_j como en (1.9). Entonces por definición (1.8) todos los elementos dentro de una racha son iguales:

$$\varepsilon_{\lambda_{j-1}^{(i)}+1} = \varepsilon_{\lambda_{j-1}^{(i)}+2} \dots = \varepsilon_{\lambda_j^{(i)}} = i \quad \begin{array}{l} j = 1, \dots, r^{(i)} \\ i = 0, 1 \end{array}$$

esto implica que

$$\sum_{k=1}^{L_j^{(1)}} \varepsilon_{\lambda_{j-1}^{(1)}+k} = L_j^{(1)}$$

y por tanto de (3.5) vale

$$(3.7) \quad \sum_{j=1}^{r^{(1)}} \sum_{k=1}^{L_j^{(1)}} \varepsilon_{\lambda_{j-1}^{(1)}+k} = \sum_{j=1}^{r^{(1)}} L_j^{(1)} = \sum_{i=1}^N \varepsilon_i = \varepsilon$$

y de aquí es válido también

$$(3.8) \quad \sum_{j=1}^{n^{(0)}} L_j^{(0)} = N - \sum_{i=1}^N \varepsilon_i = N - \varepsilon .$$

Es decir ε y $N - \varepsilon$ representan en el problema de una muestra lo que los tamaños de muestra m y n representan en el problema de dos muestras. Esto implica que las longitudes de las rachas de unos toman valores $L_j^{(1)} = 1, \dots, \varepsilon$, las longitudes de las rachas de ceros toman valores $L_j^{(0)} = 1, \dots, N - \varepsilon$ y la longitud de cualquier racha toma valores $L_j = 1, \dots, \max\{\varepsilon, N - \varepsilon\}$, $j = 1, \dots, n^{(i)}$, $i = 0, 1$

El anterior análisis tiene las siguientes consecuencias:

(i) En la construcción de una estadística para el problema de una muestra se debe utilizar ε y $N - \varepsilon$ como factores de normalización de las longitudes de las rachas.

(ii) La estadística resultante se puede utilizar también en problemas de dos muestras donde los tamaños de muestra son v.a's.

(iii) La estadística resultante es una generalización de la propuesta para el problema de dos muestras en el sentido que permite tamaños de muestra aleatorios.

A continuación se define una estadística basada en rachas para el problema de una muestra, para dos muestras pareadas y para el problema de dos muestras con tamaños de muestra aleatorios:

$$(3.9) \quad S^* = \sum_{i=1}^{n_N} \Delta_i^* h(i, L_i, n_N)$$

donde

$$\Delta_i^* = \begin{cases} \frac{1}{\varepsilon} & , \quad \varepsilon_{\lambda_{i-1}+1} = \dots = \varepsilon_{\lambda_i} = 1 \\ \frac{-1}{N-\varepsilon} & , \quad \varepsilon_{\lambda_{i-1}+1} = \dots = \varepsilon_{\lambda_i} = 0 \end{cases}$$

y la función h es como en la definición (2.6) y para el problema de una muestra se escoge siempre creciente en i y L_i .

Puesto que S^* toma sus mayores valores cuando ocurren rachas largas en los extremos de la muestra dicotomizada $\varepsilon_1, \dots, \varepsilon_N$, las pruebas basadas en S^* para probar H_0 vs. K_1 o vs. K_2 son análogas a (2.8) y (2.9).

Observación 4. Realmente Δ_i^* está definida solo en casi toda parte por causa de (3.6).

En particular eligiendo

$$h(i, L_i, n_N) = \frac{iL_i}{n_N - 1}$$

se obtiene una estadística equivalente a (2.12) para el problema de una muestra.

Observación 5. Nótese que ε en (3.7) representa la suma de las longitudes de las rachas de tipo 1 es decir de rachas de elementos de la muestra que esta por encima de la mediana.

Por otro lado ε es la estadística que se utiliza en la prueba del signo para H_0 vs. K_1 o vs. K_2 . Por lo tanto desde el punto de vista de la teoría de rachas, la prueba del signo solo contiene una mínima parte de la información necesaria para la prueba de H_0 vs. K_1 o vs. K_2 .

Muestras Pareadas.

La información disponible en este caso se puede representar por N vectores aleatorios bidimensionales independientes e idénticamente distribuidos $(X_1, Y_1), \dots, (X_N, Y_N)$. Entonces como es corriente en el caso de muestras pareadas se definen las diferencias

$$W_i = X_i - Y_i, \quad i = 1, \dots, N.$$

Cuando las v.a's. W_i , $i = 1, \dots, N$ son independientes e idénticamente distribuidas con función de distribución continua F , se procede de la misma forma que para el problema de una muestra tomando $\mu = 0$ en (3.2), (3.3) y (3.4).

Surge por otra parte, de manera natural la necesidad de una medida de la correlación entre las variables X y Y . En seguida se define un coeficiente de correlación de tipo Pearson basado en rachas.

Sean μ_x y μ_y las medianas de X y Y respectivamente. Entonces la muestra dicotomizada se puede representar por $(\epsilon_1, \epsilon_1), \dots, (\epsilon_N, \epsilon_N)$, donde

$$\epsilon_j^x = \begin{cases} 1 & \text{si } X_j > \mu_x \\ 0 & \text{si } X_j < \mu_x, \end{cases} \quad j = 1, \dots, N$$

$$\epsilon_j^y = \begin{cases} 1 & \text{si } Y_j > \mu_y \\ 0 & \text{si } Y_j < \mu_y, \end{cases} \quad j = 1, \dots, N$$

Se define $I_1 = 1$,

$$I_j^x = \begin{cases} 1, & \text{si } \varepsilon_{j-1}^x \neq \varepsilon_j^x \\ 0, & \text{si } \varepsilon_{j-1}^x = \varepsilon_j^x, \end{cases} \quad j = 1, \dots, N$$

$$I_j^y = \begin{cases} 1, & \text{si } \varepsilon_{j-1}^y \neq \varepsilon_j^y \\ 0, & \text{si } \varepsilon_{j-1}^y = \varepsilon_j^y, \end{cases} \quad j = 1, \dots, N$$

y el número de rachas de X (de Y) hasta la j -ésima observación de la muestra dicotomizada $\varepsilon_1, \dots, \varepsilon_N$ sin tener en cuenta el número de rachas de Y (de X) se define por

$$r_j^x = \sum_{k=1}^j I_k^x, \quad j = 1, \dots, N$$

$$r_j^y = \sum_{k=1}^j I_k^y, \quad j = 1, \dots, N.$$

Entonces cuando X y Y estén correlacionadas positivamente (negativamente) un aumento en r_j^x implicará un aumento (disminución) en r_j^y , $j = 1, \dots, N$. Por esta razón el coeficiente

$$R_{x,y} = \frac{\sum_{j=1}^N \left(r_j^x - \frac{r_{N+1}^x}{2} \right) \left(r_j^y - \frac{r_{N+1}^y}{2} \right)}{\left[\sum_{j=1}^N \left(r_j^x - \frac{r_{N+1}^x}{2} \right)^2 \sum_{j=1}^N \left(r_j^y - \frac{r_{N+1}^y}{2} \right)^2 \right]^{1/2}}$$

es una medida del grado de correlación lineal entre X y Y .

4. Problema de K-muestras.

Para el planteo de este problema es necesario introducir alguna notación adicional.

Sean

$$K, m_k \in \mathbb{N}, \quad k = 1, \dots, K, \quad N_0 = m_0 = 0,$$

$$(4.1) \quad N_k = \sum_{j=1}^k m_j, \quad k = 1, \dots, K, \quad N_k = \tilde{N},$$

obviamente $K < \tilde{N}$.

Las K muestras o muestra combinada se representan en este caso por las v.a's. independientes $X_{1,1}, \dots, X_{1,N_1}, X_{2,N_1+1}, \dots, X_{2,N_2}, \dots, X_{K,N_{k-1}+1}, \dots, X_{K,\tilde{N}}$ sobre un espacio de probabilidad (Ω, \mathcal{A}, P) , donde la k -ésima muestra $X_{k,N_{k-1}+1}, \dots, X_{k,N_k}$ es de tamaño $m_k = N_k - N_{k-1}$ y tiene función de distribución continua F_k , $k = 1, \dots, K$. Las estadísticas de orden se denotan por $X_{1;\tilde{N}}, \dots, X_{\tilde{N};\tilde{N}}$ y los antirangos por $D_1, \dots, D_{\tilde{N}}$.

Sea $\Omega_B = \{B_1, \dots, B_K\}$ un conjunto de K símbolos diferentes con los cuales se representan las K muestras de la siguiente manera:

$$(4.2) \quad \tilde{\eta}_j = B_k, \quad \text{si } N_{k-1}+1 \leq D_k \leq N_k, \\ k = 1, \dots, K, \quad j = 1, \dots, \tilde{N}$$

Es decir $\tilde{\eta}_j$, $j = 1, \dots, \tilde{N}$ es el símbolo B_k si la j -ésima esta

dística de orden viene de la j -ésima muestra $K = 1, \dots, K$. A la sucesión $\tilde{\eta}_1, \dots, \tilde{\eta}_N$ se le llama la muestra policotomizada puesto que representa cada una de las K muestras a través de un símbolo diferente.

Por medio de las v.a'a. $\tilde{I}_i = 1$ e

$$(4.3) \quad \tilde{I}_j = \begin{cases} 1 & \tilde{\eta}_{j-1} \neq \tilde{\eta}_j \\ 0 & \tilde{\eta}_{j-1} = \tilde{\eta}_j \end{cases} \quad j = 2, \dots, N$$

se define $\tilde{\kappa}_j$ el número de rachas hasta el j -ésimo elemento de la muestra policotomizada por

$$(4.4) \quad \begin{aligned} \tilde{\kappa}_j &= \sum_{t=1}^j \tilde{I}_t & j = 1, \dots, \tilde{N} \\ &= 1 + \sum_{t=2}^j \tilde{I}_t & j = 2, \dots, \tilde{N} . \end{aligned}$$

En particular cuando $K = 2$ se obtiene κ_j , $j = 1, \dots, N$ de (1.7).

Observación 6. Para $K > 2$, $\tilde{\kappa}_N$: el número total de rachas en la muestra policotomizada $\tilde{\eta}_1, \dots, \tilde{\eta}_N$, es la generalización para K muestras de la estadística propuesta por Wald y Wolfowitz (1940). Barton y David (1959), obtienen también $\tilde{\kappa}_N$ por un método diferente y dan una aproximación de su distribución a través de la distribución de Poisson.

Observación 7. $\tilde{\eta}_j$, \tilde{I}_j y $\tilde{\kappa}_j$, $j = 1, \dots, \tilde{N}$ se pueden obtener también de manera natural desde los rasgos de las observaciones de la muestra combinada de la misma forma que η_j , I_j y κ_j en (1.5), (1.6) y (1.7).

En este párrafo se consideran las siguientes hipótesis:

Hipótesis nula.

$$(4.5) \quad H_0 : F_1(x) = \dots = F_K(x), \quad \forall x \in \mathbb{R}.$$

Alternativa general.

$$(4.6) \quad \tilde{K}_1 : F_i(x) \neq F_j(x)$$

para algún par i, j , $1 \leq i \neq j \leq K$ y para algún $x \in \mathbb{R}$.

Alternativa de localización de dos colas.

$$(4.7) \quad \tilde{K}_2 : F_i(x) = F_j(x - \mu)$$

para algún par i, j , $1 \leq i \neq j \leq K$, $\mu \neq 0$, $\forall x \in \mathbb{R}$.

Alternativa de escala de dos colas

$$(4.8) \quad \tilde{K}_3 : F_i(x) = F_j\left(\frac{x - \theta}{\zeta}\right)$$

para algún par i, j , $1 \leq i \neq j \leq K$, $\zeta \neq 0$, $\zeta \neq 1$, $\forall x \in \mathbb{R}$, donde $\theta = \mu_j - \mu_i$, $1 \leq i \neq j \leq K$ y μ_k es la mediana de F_k , $k = 1, \dots, K$. Esta manera de tomar θ significa que F_i y F_j deben tener la misma mediana.

La estadística general para el problema de K muestras se define por

$$(4.9) \quad \tilde{T}^* = \sum_{i=1}^{\tilde{n}_N} \tilde{\delta}_i^* h(i, i, \tilde{n}_N),$$

donde $\tilde{\tau}_N$ es como en (4.4), L_i , $i = 1, \dots, \tilde{\tau}_N$ representa la longitud de la i -ésima racha en la muestra policotomizada $\tilde{\eta}_1, \dots, \tilde{\eta}_{\tilde{N}}$, $L_i = 1, \dots, \max\{m_1, \dots, m_k\}$ y $\tilde{\delta}_i^*$ se define a continuación

$$(4.10) \quad \tilde{\delta}_i^* = \frac{1}{m_k} \quad \text{si} \quad \tilde{\eta}_{\lambda_{i-1}+1} = \dots = \tilde{\eta}_{\lambda_i} = B_k,$$

$$k = 1, \dots, K, \quad i = 1, \dots, \tilde{\tau}_N$$

La función h es similar a la función h utilizada para T^* en (2.6) y se escoge también análogamente para las alternativas de localización y escala mientras que para la alternativa general \tilde{K}_1 se debe tomar

$$(4.11) \quad h(i, L_i, \tilde{\tau}_N) = m_i, \quad i = 1, \dots, \tilde{\tau}_N$$

lo cual implica $\tilde{T}^* = \tilde{\tau}_N$.

(i) Pruebas para la Alternativa General \tilde{K}_1 :

Bajo H_0 se espera que haya muchas rachas en la muestra dicotomizada. Por lo tanto valores pequeños de $\tilde{T}^* = \tilde{\tau}_N$ apoyarán la alternativa \tilde{K}_1 lo cual se puede interpretar como una prueba aleatorizada de nivel $\alpha \in (0, 1)$ para H_0 vs. \tilde{K}_1 :

$$(4.12) \quad \tilde{\phi}_1 = \begin{cases} 1, & \text{si } \tilde{\tau}_N^* < c \\ \gamma, & \text{si } \tilde{\tau}_N^* = c \\ 0, & \text{si } \tilde{\tau}_N^* > c \end{cases}$$

donde γ y c se determinan de manera análoga a los de (2.8).

(ii) Pruebas para la Alternativa de Localización de dos colas.

En este caso la función h se escoge creciente en i y L_i .

Por esta razón \tilde{T}^* toma sus mayores valores cuando hay rachas grandes en los extremos de la muestra dicotomizada. Entonces se tiene una prueba aleatorizada de nivel $\alpha \in (0,1)$ para H_0 vs. K_2 :

$$(4.13) \quad \tilde{\phi}_2 = \begin{cases} 1 & \text{si } \tilde{T}^* > c \\ \gamma & \text{si } \tilde{T}^* = c \\ 0 & \text{si } \tilde{T}^* < c \end{cases}$$

donde γ y c se determinan de manera análoga a los de (2.8).

(iii) Pruebas para la Alternativa de dos colas.

En este caso la función h se escoge creciente en $\left| i - \frac{N+1}{2} \right|$ y L_i . Por esta razón \tilde{T}^* toma sus mayores valores cuando hay rachas grandes en los extremos de la muestra dicotomizada. Entonces se tiene una prueba aleatorizada de nivel $\alpha \in (0,1)$ para H_0 vs. K_3 :

$$(4.14) \quad \tilde{\phi}_2 = \begin{cases} 1 & \text{si } \tilde{T}^* > c \\ \gamma & \text{si } \tilde{T}^* = c \\ 0 & \text{si } \tilde{T}^* < c \end{cases}$$

donde γ y c se determinan de manera análoga a los de (2.8).

5. Pruebas Insesgadas.

En este parágrafo se demuestra que una subfamilia de las pruebas aquí presentadas son insesgadas. Para esto se hace primero una transformación de las sumas aleatorias (2.6), (3.9) y (4.9) utilizadas en las pruebas, en sumas con núme-

ro fijo de sumandos, resultado que se obtiene como aplicación directa de la transformación (1.14).

Sea

$$(5.1) \quad Q^* = \sum_{i=1}^{r_N} \omega_i^* g(i, r_N) L_i$$

la forma general de cualquiera de las estadísticas T^* , S^* o \tilde{T}^* cuando se escoge la función h de la forma

$$(5.2) \quad h(i, L_i, r_N) = g(i, r_N) L_i. \quad i = 1, \dots, r_N,$$

donde g es una función real no negativa, acotada y se escoge creciente en i para problemas de localización en tanto que para problemas de escala se toma creciente en $\left| i - \frac{r_N + 1}{2} \right|$, y los valores de ω_i^* se determinan adecuadamente para T^* , S^* o \tilde{T}^* como ya se definieron en (2.6), (3.9) y (4.10).

Se define la estadística

$$(5.3) \quad T = \sum_{i=1}^N \delta_i g(r_i, r_N)$$

donde r_j , $j = 1, \dots, N$ es como en (1.7) para problemas de dos y una muestra o su equivalente (4.4) para el caso de K muestras. Además la función δ_i se escoge de la siguiente manera:

Para el problema de dos muestras

$$(5.4) \quad \delta_j = \begin{cases} \frac{1}{m}, & \text{si } \eta_j = 1 \\ -\frac{1}{m}, & \text{si } \eta_j = 0, \quad j = 1, \dots, N \end{cases}$$

Para el problema de una muestra

$$(5.5) \quad \delta_j = \begin{cases} \frac{1}{\varepsilon}, & \text{si } \varepsilon_j = 1 \\ -\frac{1}{N-\varepsilon}, & \text{si } \varepsilon_j \neq 1, \end{cases} \quad j = 1, \dots, N$$

Para el problema de K muestras

$$(5.6) \quad \delta_j = \frac{1}{m_k} \quad \text{si } \tilde{\eta}_j = B_k, \quad k = 1, \dots, K, \quad j = 1, \dots, \tilde{N}$$

de acuerdo con las definiciones el valor de δ_j es el mismo para todos los elementos dentro de una racha:

$$\delta_j = k \iff \delta_{\lambda_{j-1}+1} = \dots = \delta_{\lambda_j} = k, \quad j = 1, \dots, N.$$

Entonces de (1.13)

$$\omega_i^* g(i, n_N) L_i = \sum_{k=1}^{L_i} \delta_{\lambda_{j-1}+k} g(n_{\lambda_{j-1}+k}, n_N)$$

y de (1.14)

$$(5.7) \quad \sum_{i=1}^{n_N} \omega_i^* g(i, n_N) L_i = \sum_{i=1}^N \delta_i g(n_i, n_N)$$

es decir $Q^* = T$.

Observación 8. La demostración es igualmente válida si en (5.4) y (5.5) se toman valores positivos $1/n$ y $1/(N-\varepsilon)$ para δ_i , $i = 1, \dots, N$ cuando $\eta = 0$.

La propiedad de insesgamiento de las pruebas presentadas aquí se demostrará para hipótesis más generales que las

presentadas en los párrafos 2, 3 y 4. Para su formulación se requiere la siguiente definición.

DEFINICION. Sean F y G las funciones de distribución de las v.a's. X y Y respectivamente. Se dice que F es estocásticamente mayor o igual que G (o bien que X es estocásticamente mayor o igual que Y) y se escribe

$$(5.8) \quad F \cdot \geq G \quad (\text{o bien } X \cdot \geq Y)$$

si y solo si

$$F(x) \leq G(x) \quad \forall x \in \mathbb{R}.$$

Análogamente se definen las relaciones $F \cdot > G$, $F \leq \cdot G$ y $F < \cdot G$.

En este párrafo se consideraran las siguientes hipótesis:

Hipótesis nula.

$$(5.9) \quad H_0^* : F \cdot \geq G \quad (\text{o bien } X \cdot \geq Y).$$

Hipótesis alternativa

$$(5.10) \quad K : F < \cdot G.$$

En el siguiente teorema se demuestra que a partir de la familia de estadísticas

$$(5.11) \quad S = \sum_{i=1}^N b_i g(r_{R_i}, r_N)$$

y r_{R_i} es como se definió en (1.7), se obtiene una familia de pruebas insesgadas para H_0^* en (5.9).

Las constantes de regresión en (5.11) b_i , $i = 1, \dots, N$ se escogen de acuerdo con el problema, así:

(i) Para el problema de dos muestras:

$$b_i = \begin{cases} \frac{1}{m}, & \text{si } i = 1, \dots, m \\ \frac{1}{n}, & \text{si } i = m+1, \dots, N, \quad i = 1, \dots, N. \end{cases}$$

(ii) Para el problema de una muestra:

$$b_i = \begin{cases} \frac{1}{\varepsilon}, & \text{si } i = 1, \dots, \\ \frac{1}{N-\varepsilon}, & \text{si } i = \varepsilon+1, \dots, N, \quad i = 1, \dots, N \end{cases}$$

(iii) Para el problema de K -muestras:

$$b_i = \frac{1}{m_k} \quad \text{si } N_{k-1}+1 \leq i \leq N_k, \quad k = 1, \dots, K, \quad i = 1, \dots, \tilde{N}$$

TEOREMA. La Prueba basada en la estadística S de (5.11):

$$(5.13) \quad \phi = \begin{cases} 1, & \text{para } S > C_\alpha \\ \gamma, & \text{para } S = C_\alpha \\ 0, & \text{para } S < C_\alpha \end{cases}$$

donde C_α se determina de tal manera que

$$(5.14) \quad E_{H_0^*}(\phi) = \alpha,$$

es una prueba insesgada de nivel $\alpha \in (0,1)$ para H_0^* versus K .

Antes de hacer la demostración se presentan dos lemas que la hacen más corta.

LEMA 1. Sean X y Y v.a's con distribuciones P_x y P_y respectivamente y sea la función $h: \mathbb{R} \rightarrow \mathbb{R}$ monótona creciente.

Si

$$P_x \leq P_y \quad \text{entonces} \quad P_{h(x)} \leq P_{h(y)}$$

LEMA 2. Sean X y Y como en el Lema 1. Si

$$P_x \leq P_y \quad \text{entonces} \quad E(X) \leq E(Y)$$

siempre que $E(X)$ y $E(Y)$ sean finitas.

Demostración.

(i) Problema de dos muestras.

Sean F y G las funciones de distribución continuas de X_1, \dots, X_m , y de X_{m+1}, \dots, X_N respectivamente.

Sean U_1, \dots, U_m , y U_{m+1}, \dots, U_N v.a's. independientes e idénticamente distribuidas con función de distribución uniforme en $(0,1)$. Entonces X_1, \dots, X_m , y X_{m+1}, \dots, X_N tienen la misma distribución que $F^{-1}(U_i)$, $i = 1, \dots, m$ y que $G^{-1}(U_i)$, $i = m+1, \dots, N$ respectivamente.

Se definen las v.a's:

$$V_i = F^{-1}(U_i), \quad i = 1, \dots, N$$

y

$$w_i = \begin{cases} F^{-1}(u_i) & i = 1, \dots, m \\ G^{-1}(u_i) & i = m+1, \dots, N. \end{cases}$$

Entonces para $V = (V_1, \dots, V_N)$ valen las siguientes relaciones:

$$\begin{aligned} P_x &= P(V_1, \dots, V_m, V_{m+1}, \dots, V_N) \\ P_x^{H_0} &= P_V^{H_0} \\ P_{R(X)}^{H_0} &= P_{R(V)}^{H_0} \end{aligned}$$

donde P^{H_0} es la distribución bajo $H_0 : F = G$ del vector indicado en el subíndice.

Si $F < G$ entonces por definición $F(x) \geq G(x)$ y por tanto $F^{-1}(y) \leq G^{-1}(y)$ para todo $y \in (0,1)$ y esto implica

$$\begin{aligned} V_i &= W_i, & i = 1, \dots, m \\ V_i &\leq W_i, & i = m+1, \dots, N. \end{aligned}$$

Por otra parte para todo $i, j > m$, $W_j \geq W_i$ si y sólo si $U_j \geq U_i$ o sea si $V_j \geq V_i$. De aquí se deduce, para $i > m$, que

$$\begin{aligned} R_i(V) &= \sum_{k=1}^N U(X_i - X_k) = \sum_{k=1}^m U(X_i - X_k) + \sum_{k=m+1}^N U(X_i - X_k) \\ &\leq \sum_{k=1}^m U(W_i - W_k) + \sum_{k=m+1}^N U(W_i - W_k) = R_i(W) \end{aligned}$$

donde $\underline{w} = (w_1, \dots, w_N)$ y $U(x) = 1$ ó 0 según $x \geq 0$ ó $x < 0$.

Por (1.4), r_{R_i} es una función no decreciente de R_i por lo tanto para $i > m$ se tiene:

$$r_{R_i}(V) \leq r_{R_i}(W) \quad (\text{c.s.: casi seguro})$$

Por lo tanto para g como en (5.2) y $j > m$

$$(5.15) \quad g(r_{R_i}(V), r_N) \leq g(r_{R_i}(W), r_N)$$

$$(5.16) \quad \sum_{i=1}^N b_i g(r_{R_i}(V), r_N) \leq \sum_{i=1}^N b_i g(r_{R_i}(W), r_N)$$

Por tanto bajo $H_0 : F = G$ y K de (5.10) y por Lema 1 se cumple:

$$(5.17) \quad P_S^{H_0} \leq P_S^k$$

donde P_V^k es la distribución de S bajo la alternativa K . Además bajo H_0^* vale también

$$P_S^{H_0} \geq P_S^{H_0^*}$$

entonces por el Lema 2:

$$P_{\phi(S)}^{H_0} \geq P_{\phi(S)}^{H_0^*}$$

Así que por (5.14) y el Lema 2 se tiene

$$(5.18) \quad E_{H_0^*} \phi(S) < \alpha.$$

Es decir de (5.14), Lema 1 y Lema 2 sigue análogamente:

$$(5.19) \quad E_k \phi(S) \geq \alpha .$$

(ii) Problema de una Muestra.

La demostración es exactamente igual hasta la ecuación (5.15) y de ahí en adelante las ecuaciones valen CS.

(iii) Problema de K-Muestras.

Análogo al problema de una muestra.

Conclusiones.

Las subfamilias de pruebas de rachas aquí propuestas tienen las siguientes características generales:

(i) Son una alternativa que como lo muestran los casos especiales trabajados en Ortiz (1983), Corzo-Ortiz (1983), Fernández-Ortiz (1986), tiene muy buen comportamiento en términos de su potencia y de su convergencia a la distribución normal, frente a las pruebas basadas en rangos. Esto hace pensar que si en otras aplicaciones se mantiene este comportamiento, lo que se ha descubierto es una poderosa familia de pruebas basadas en rachas para alternativas de localización y escala para cada uno de los problemas analizados.

(ii) Aunque no se mantuviera el comportamiento analizado en los trabajos arriba mencionados, las familias propuestas constituyen la única alternativa disponible para problemas de localización y escala en situaciones en que se dispone de información que solo se encuentra en escala nominal.

BIBLIOGRAFIA

- Barton, D.E., David, F.N. (1957). "Multiple Runs". *Biometrika*, 44, 168-178.
- Corzo, J. Ortiz, J. (1983). "Una prueba de dispersión basada en secuencias". *Revista Colombiana de Estadística* 8, 34-38.
- Corzo, J. (1989). "Verallgemeinerte Runtests für Lage - und Scalen Alternativen". Tesis Doctoral, Univ. Dortmund. R.F.A.
- Corzo, J. (1989). "Teoría de Rachas". *Revista Colombiana de Estadística*. N° 19-20.
- Fernández, Ortiz, J. (1986). "Pruebas de Hipótesis basadas en Secuencias". *Revista Colombiana de Estadística* N°12.
- Ortiz, J. (1983). "Pruebas de hipótesis sobre parámetros de localización basadas en secuencias". *Revista Colombiana de Estadística* N° 8.

*