

STAT 824 sp 2023 Lec 12 slides

Wilcoxon rank-sum test

Karl B. Gregory

University of South Carolina

These slides are an instructional aid; their sole purpose is to display, during the lecture, definitions, plots, results, etc. which take too much time to write by hand on the blackboard. They are not intended to explain or expound on any material.

Table of Contents

1 Wilcoxon rank sum test

2 Power comparisons

Suppose we collect random samples from “control” and “treatment” populations:

$$X_1, \dots, X_m \stackrel{\text{ind}}{\sim} F \quad \text{“control”}$$

$$Y_1, \dots, Y_n \stackrel{\text{ind}}{\sim} G \quad \text{“treatment”}$$

We wish to test for treatment effectiveness (are Y’s bigger than X’s?).

Wilcoxon rank sum test (quintessential nonparametric test)

The *Wilcoxon rank sum test (WXRS)* concludes a “positive treatment effect” if

$$W_{XY} = \sum_{i=1}^m \sum_{j=1}^n \mathbf{1}(X_i < Y_j) \geq c,$$

where c can be calibrated to control the Type I error rate.

Can modify to find a “negative” or “either direction” treatment effect.

If $G = F$, the (null) distribution of W_{XY} is the same for any continuous F .

Rank-sum form of Wilcoxon rank sum statistic

An alternate way of computing W_{XY} :

- 1 Sort all the data $(X_1, \dots, X_m, Y_1, \dots, Y_n)$
- 2 Obtain the ranks.
- 3 Keep the ranks corresponding to Y_1, \dots, Y_n , calling these S_1, \dots, S_n .

Then $W_{XY} = S_1 + \dots + S_n - n(n+1)/2$.

Let $W_S = S_1 + \dots + S_n$.

Exercise: Show that $W_{XY} = W_S - n(n+1)/2$.

Theorem

Let $X_1, \dots, X_m, Y_1, \dots, Y_n$ be continuous iid rvs and set $N = n + m$. Then

$$P(\{S_1, \dots, S_n\} = \{s_1, \dots, s_n\}) = \frac{1}{\binom{N}{n}}$$

for all sets of n ranks $\{s_1, \dots, s_n\} \subset \{1, \dots, N\}$.

Exercise: Tabulate the null distribution of W_{XY} under $N = 5, n = 2$.

Example code:

```

m <- 25
n <- 20
X <- rnorm(m,1,1)
Y <- rnorm(n,1,1)
N <- n + m

U <- c(X,Y)
id <- c(rep(1,m),rep(2,n))
id_ord <- id[order(U)]
S <- c(1:N)[which(id_ord == 2)]
Ws <- sum(S)
Wxy <- Ws - n*(n+1)/2

# note that n and m are interchangeable:
1 - pwilcox(Wxy-1,m = n,n = m) # these are slow for large n, m
1 - pwilcox(Wxy-1,m = m,n = n)

EWs <- n*(N+1)/2
VarWs <- n*(N-n)*(N+1)/12
1 - pnorm(Ws,EWs,sqrt(VarWs))
1 - pnorm(Ws,EWs + 1/2,sqrt(VarWs)) # with continuity correction

wilcox.test(x = Y, y = X, alternative = "greater")

```

On computing the exact distribution of W_{XY} when N and n are large. . .



Theorem (Asymptotic Normality of rank sum under the null)

Under $H_0: F = G$ we have

$$\frac{W_S - \mathbb{E}W_S}{\sqrt{\text{Var } W_S}} \xrightarrow{D} \text{Normal}(0, 1)$$

as $N \rightarrow \infty$, provided $n \rightarrow \infty$ and $N - n \rightarrow \infty$.

Exercise: Show $\mathbb{E}W_S = \frac{1}{2}n(N + 1)$ and $\text{Var } W_S = \frac{1}{12}n(N - n)(N + 1)$.

Corollary

An asymptotic p -value for testing $H_0: F = G$ versus the “right-sided” alternative is

$$1 - \Phi \left(\frac{W_S - \frac{1}{2}n(N+1) + \frac{1}{2}}{\sqrt{\frac{1}{12}n(N-n)(N+1)}} \right),$$

where the extra $\frac{1}{2}$ is a “continuity correction”.

Sketch of asymptotic Normality proof

Assume $H_0: F = G$ and introduce $U_1, \dots, U_n \stackrel{\text{ind}}{\sim} \text{Uniform}(0, 1)$. Then:

- 1 Write W_S as a sum of *dependent* rvs: $W_S = \sum_{i=1}^N i \cdot J_i$, $J_i = \mathbf{1}(U_i \leq U_{(n)})$.
- 2 Introduce approximator \tilde{W}_S , which is a sum of *independent* rvs:

$$\tilde{W}_S = \sum_{i=1}^N \left(i - \frac{N+1}{2}\right) K_i + \frac{n(N+1)}{2}, \quad K_i = \mathbf{1}(U_i \leq n/N).$$

- 3 Show that $\frac{\tilde{W}_S - \mathbb{E}\tilde{W}_S}{\sqrt{\text{Var } \tilde{W}_S}} \xrightarrow{D} \text{Normal}(0, 1)$ as $N, n, m \rightarrow \infty$.
- 4 Argue same holds for W_S since $\frac{\mathbb{E}(\tilde{W}_S - W_S)^2}{\text{Var } \tilde{W}_S} \rightarrow 0$ as $N, n, m \rightarrow \infty$.

Exercise:

- 1 Show $\mathbb{E}\tilde{W}_S = \mathbb{E}W_S$ and $\text{Var } \tilde{W}_S = \frac{(N-1)}{N} \text{Var } W_S$.
- 2 Show $\frac{\tilde{W}_S - \mathbb{E}\tilde{W}_S}{\sqrt{\text{Var } \tilde{W}_S}} \xrightarrow{D} \text{Normal}(0, 1)$ as $N, n, m \rightarrow \infty$ with Lindeberg CLT.

1 Wilcoxon rank sum test

2 Power comparisons

To analyze the power of the WXRS we must specify an alternative to $H_0: F = G$.

Location shift model

In the *location-shift* model we assume $G(x) = F(x - \Delta)$ for some $\Delta \in \mathbb{R}$.

We will consider the right-sided test $H_0: \Delta \leq 0$ vs $H_1: \Delta > 0$.

Exercise: Show that the power of the rule $W_{XY} \geq c$ is nondecreasing in Δ .

It is convenient to use a Normal approximation to the power:

Theorem (Approximate power of WXRS in location-shift model)

In the location-shift model the power of $W_{XY} \geq c$ admits the approximation

$$\gamma(\Delta) \approx 1 - \Phi \left(\frac{c - nmp_1(\Delta)}{\sqrt{\vartheta(\Delta)}} \right)$$

provided N , n , and $N - n$ are all large, where $p_1(\Delta) = P(X_1 < Y_1)$ and

$$\vartheta(\Delta) = mnp_1(\Delta)[1 - p_1(\Delta)] + mn(n-1)[p_2(\Delta) - p_1^2(\Delta)] + nm(m-1)[p_3(\Delta) - p_1^2(\Delta)]$$

with $p_2(\Delta) = P(X_1 < Y_1, X_2 < Y_1)$ and $p_3(\Delta) = P(X_1 < Y_1, X_1 < Y_2)$.

Exercise:

- 1 Establish the above result.
- 2 Find the value of c such that the test has size approximately equal to α .

Exercise: Show that making the substitutions

- 1 $c = c_\alpha$
- 2 $p_1(\Delta) = 1/2 + \Delta f^*(0)$, f^* the density of $X_1 - X_2$
- 3 $\vartheta(\Delta) = \vartheta(0)$

leads to the approximate power curve for the size- α test given by

$$\tilde{\gamma}_\alpha(\Delta) = 1 - \Phi \left(z_\alpha - \sqrt{\frac{12nm}{N+1}} \cdot \Delta \cdot f^*(0) \right).$$

Exercise:

- 1 Show that if F is Normal, $n = m$, and $N + 1$ is replaced by $2n$, we obtain

$$\tilde{\gamma}_\alpha(\Delta) = 1 - \Phi \left(z_\alpha - \frac{\sqrt{6n}}{2\sigma\sqrt{\pi}} \cdot \Delta \right).$$

- 2 Find the smallest n such that the WXRS has power $\geq \gamma^*$ for all $\Delta \geq \Delta^*$.
- 3 Compare to n needed for the equal-variances two-sample z -test.

