

Hypothesis Testing 4

Day 14 (2/20/20)

The Wilcoxon Rank-Sum test

Suppose we have two independent random samples X_1, \dots, X_m from population 1 with continuous cdf $F(\cdot)$ and Y_1, \dots, Y_n from population 2 with continuous cdf $G(\cdot)$. Assume that the distributions are the same except for a shift in mean, i.e.

$$Y \stackrel{d}{=} X + \Delta$$

For testing

$$H_0 : \Delta = 0 \text{ vs } H_1 : \Delta \neq 0$$

the Wilcoxon rank sum test does the following:

1. Order the combined sample of $N = m + n$ X -values and Y -values from smallest to largest. Let S_1 denote the rank of Y_1, \dots, S_n denote the rank of Y_n in this joint ordering.
2. Let W be the sum of the ranks of the Y values

$$W = \sum_{i=1}^n S_i$$

3. Large-sample test: Reject H_0 if

$$|W^*| = \left| \frac{W - [n(m+n+1)/2]}{\sqrt{mn(m+n+1)/12}} \right| \geq z_{\alpha/2}$$

4. Calculation in R

```
> attach(aztdoses)
> x1<-azt[dose==300]
> x2<-azt[dose==600]
> wilcox.test(x1,x2)
```

```
Wilcoxon rank sum test with continuity
correction
```

```
data: x1 and x2
W = 23, p-value = 0.04459
alternative hypothesis: true location shift is not equal to 0
```

```
Warning message:
In wilcox.test.default(x1, x2) : cannot compute exact p-value with ties
```

5. Manual calculation in R

```
> c(x1,x2)
[1] 284 279 289 292 287 295 285 279 306 298 298 307
[13] 297 279 291 335 299 300 306 291
> r<-rank(c(x1,x2))
> r
[1] 4.0 2.0 7.0 10.0 6.0 11.0 5.0 2.0 17.5
[10] 13.5 13.5 19.0 12.0 2.0 8.5 20.0 15.0 16.0
[19] 17.5 8.5
> sum(r[1:10])
[1] 78
> (78-10*(20+1)/2)/sqrt(10*10*(20+1)/12)
[1] -2.041008
> 2*(1-pnorm(2.041008))
[1] 0.04125003
```