

# Weighted U-statistics for Simple Linear Regression

Kim Do-hyun, Kim Chul-soo\*

단순한 선형회귀에 대한 무게있는 U-통계량

金道鉉, 金鐵洙\*

## 1. Introduction

Consider an arbitrary kernel  $h(x_1, \dots, x_n)$ , not necessarily symmetric, to be applied as usual to observations  $X_1, \dots, X_n$  taken  $m$  at a time. Suppose also that each term  $h(X_{i_1}, \dots, X_{i_m})$  becomes weighted by a factor  $w(i_1, \dots, i_m)$  depending only on the indices  $i_1, \dots, i_m$ . In this case the U-statistics sum takes the more general form  $T_n = w(i_1, \dots, i_m)h(X_{i_1}, \dots, X_{i_m})$ . In this case that  $h$  is symmetric and the weights  $w(i_1, \dots, i_m)$  take only 0 or 1 as values. Certain "weighted U-statistics" for simple linear regression take the form  $T_n$ .

Consider the simple linear regression model  $Y_i = \alpha + \beta x_i + e_i, 1 \leq i \leq n, (1.1)$

where  $\alpha$  and  $\beta$  are unknown parameters,  $x_i$  are known regression scores, and  $e_i$  are i. i. d. random variables with c. d. f.  $F$ .

In this paper we extend the procedures of Theil by using the projection of weighted U-statistic of the form weighted rank for the simple linear regression model (1.1).

## 2. A weighted U-statistics

For the regression model (1.1), assume that  $F$  is continuous to rule out ties among the  $Y$ 's. Also assume that  $x_1 < x_2 < \dots < x_n$  with at least one strict inequality. We will consider inferences for based on weighted U-statistics defined by

$$T_{\beta} = \sum_{i=1}^{n-1} \sum_{j=i+1}^n a_{ij} \phi(Y_i - \alpha - \beta x_i, Y_j - \alpha - \beta x_j),$$

where  $\phi(u, v) = 1$  or  $0$  according as  $u < v$  or  $u > v$ . The weights  $a_{ij} > 0$  are arbitrary but assume that  $a_{ij} = 0$  if  $x_i = x_j$ . We define the slope of the line segment from the point  $(x_i, Y_i)$  to the point  $(x_j, Y_j)$  by

$$S_{ij} = (Y_j - Y_i) / (x_j - x_i), 1 \leq j, x_i \neq x_j.$$

Note: 1.  $T_{\beta}$  is a function of the slope  $S_{ij}$  since  $\phi(Y_i - \alpha - \beta x_i, Y_j - \alpha - \beta x_j) = 1$  when  $S_{ij} > \beta$ .

2. The distribution of  $T_{\beta}$  depends on the weights  $a_{ij}$ .
3. If  $a_{ij} = 1$ , the wilcoxon distribution can be applied.

But it is not generally feasible to tabulate the exact distribution for smaller sample sizes. For larger sample sizes, the distribution of  $T_{\beta}$  can be approximated by a normal distribution.

### 2.1. An estimator associated with the Theil

-----  
 師範大學 專任講師, 延世大學校 講師 \*

statistic to estimate  $\beta$  of model (1.1).

a) Let  $N = \binom{n}{2}$  and form the  $N$  sample slope

$$S_{ij} = (Y_j - Y_i) / (x_j - x_i), \quad i < j, \quad x_i \neq x_j$$

b) The estimator of  $\beta$  is

$$\hat{\beta} = \text{median } \{S_{ij}\}.$$

Let  $S^{(1)} < S^{(2)} < \dots < S^{(N)}$  denote the ordered values of  $S_{ij}$ .

Then if  $N$  is odd, say  $N = 2k + 1$ , we have  $\hat{\beta} = S^{(k+1)}$ .

If  $N$  is even, say  $N = 2k$ , then

$$\hat{\beta} = \frac{S^{(k)} + S^{(k+1)}}{2}$$

The estimator  $\hat{\beta}$  is less sensitive to gross errors than is the classical least-square estimator.

Sen has generalized Theil's estimator to the case where the  $x$ 's are not distinct.

### 2.2. Estimation for $\beta$ by $T_\beta$

Consider the estimator  $\hat{\beta}$  of the slope parameter defined by

$$\hat{\beta}_U = \sup \{ \beta : T_\beta \geq \frac{a_{\dots}}{2} \},$$

$$\hat{\beta}_L = \inf \{ \beta : T_\beta \leq \frac{a_{\dots}}{2} \},$$

and  $\hat{\beta} = (\hat{\beta}_U + \hat{\beta}_L) / 2$ , where

$$a_{\dots} = \sum_{i=1}^{n-1} \sum_{j=i+1}^n a_{ij}.$$

Note.  $T_\beta$  is a nonincreasing, left-continuous function of that ranges from  $a_{\dots}$  down to zero. It is a step function with jumps  $a_{ij}$  occurring at the point  $S_{ij}$ .

Consider the asymptotic distribution of  $T_\beta$ . The distribution of  $T_\beta$  when the slope parameter is  $\beta$  is the same as the distribution of  $T_0$  when  $\beta = 0$ . Now  $T_0$  is a function of the ranks of the  $Y$ 's and the basic approach is to consider the projection of

$T_0$ , say  $T_0^*$ , into the family of linear rank statistics and to establish that  $T_0$  and  $T_0^*$  have the same asymptotic distributions.

Define the row and column sums of the weights by

$$a_{i\cdot} = \sum_{j=1}^n a_{ij}$$

and

$$a_{\cdot j} = \sum_{i=1}^{j-1} a_{ij} \text{ for } 1 < i, j < n$$

with  $a_{n\cdot} = 0$  and  $a_{\cdot 1} = 0$ . Let  $A_i = a_{i\cdot} - a_{\cdot i}$  and

$$a_{\dots} = \sum_{i=1}^{n-1} \sum_{j=i+1}^n a_{ij}.$$

From the Hájek and Šidák, the projection of  $T_0$  is

$$T_0^* = (1/n) \sum_{i=1}^n A_i R_i + a_{\dots} / 2,$$

where  $R_i$  is the rank of  $Y_i$  among  $\{Y_1, \dots, Y_n\}$ ,  $1 < i < n$ . That is, this can be verified by the following lemma and theorem.

**Lemma (Hájek).** Let  $Z_1, \dots, Z_n$  be independent random variables and  $S = S(Z_1, \dots, Z_n)$  any statistic satisfying  $E(S^2) < \infty$ . Then the random variable

$$\hat{S} = \sum_{i=1}^n E(S | Z_i) - (n-1)E(S)$$

satisfies

$$E(\hat{S}) = E(S)$$

and

$$E(\hat{S} - S)^2 = \text{Var}(S) - \text{Var}(\hat{S}).$$

The random variable  $\hat{S}$  is called the projection of  $S$ .

It is also possible to apply the technique to project a statistic onto dependent random variables.

**Theorem.** Consider a rank statistic  $T = T(R_1, \dots, R_N)$  and put

$$\hat{a}(i, j) = E(T | R_i = j), \quad 1 < i, j < N.$$

Then the statistic

$$\hat{T} = \frac{N-1}{N} \sum_{i=1}^N \hat{a}(i, R_i) - (N-2)E(T)$$

is the projection of  $T$  into the family of linear rank statistics.

*Proof.* See Hájek, p. 59.

The following two theorems are immediate from

Hájek and Šidák(1967, p. 163).

Theorem 1: Assume  $\beta=0$  and the condition

$$A: \sum_{i=1}^m A_i^2 / \max_{1 \leq i \leq m} A_i^2 \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Then  $T_0^*$  is asymptotically  $N(a \dots / 2, \sum_{i=1}^n A_i^2 / 12)$ .

Theorem 2: Assume  $\beta=0$ , condition A, and condition

$$B: \sum_{i=1}^n a_{ij}^2 / \sum_{i=1}^n A_i^2 \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Then  $T_0$  is asymptotically  $N(a \dots / 2, \sum_{i=1}^n A_i^2 / 12)$ .

The exact variance of  $T_0$  is

$$(\sum_{i=1}^n A_i^2 + \sum_{i=1}^{n-1} \sum_{j=i+1}^n a_{ij}^2) / 12.$$

Now the standard test of  $\beta=0$  of the model (1.1) is introduced. This is based on  $\sum_{i=1}^n (x_i - \bar{x}) Y_i$ .

This suggests a rank statistic

$$U = \sum_{i=1}^n (x_i - \bar{x}) R_i,$$

where  $R_1, \dots, R_n$  are the ranks of  $Y_1, \dots, Y_n$ .

Under  $\beta=0$ ,  $Y_1, \dots, Y_n$  are i. i. d. random variables with c. d. f.  $F(y - a)$ , we obtain the following theorem.

Theorem 3: In the simple regression model (1.1), suppose  $\beta=0$  and suppose that

$$(1/n) \sum_{i=1}^n (x_i - \bar{x})^2 \rightarrow \delta^2 > 0$$

Then

$$U^* = \frac{1}{(n+1)\sqrt{n}} \sum_{i=1}^n (x_i - \bar{x})(R_i - (n+1)/2) \text{ is asymptotically } N(0, \delta^2/12).$$

Proof. First, note that the rank of  $Y_i$  among  $Y_1, \dots, Y_n$  can be written as

$$R_i = 1 + \sum_{j=1}^n s(Y_j - Y_i).$$

where  $s(x) = 1$  if  $x > 0$  and 0 otherwise. Then, since

$$E[s(Y_j - Y_i) | Y_k = y] = \begin{cases} F(y) & k=j \\ 1-F(y) & k=i \\ 1/2 & k \neq i \text{ or } j. \end{cases}$$

we have

$$E[R_j | Y_k = y] = 1 + \sum_{i=1}^n E[s(Y_i - Y_j) | Y_k = y] = \begin{cases} 1 + (n-1)F(y) & k=j \\ 1 + (n-2)/2 + (1-F(y)) & k \neq j \end{cases}$$

And we have

$$E[U^* | Y_k = y] = \frac{1}{(n+1)\sqrt{n}} \sum_{i=1}^n (x_i - \bar{x}) [1/2 - F(y)] + (x_k - \bar{x}) [(n-1)F(y) - (n-1)/2] = \frac{\sqrt{n}}{(n+1)} (x_k - \bar{x}) [F(y) - 1/2].$$

Hence the projection  $V_p$  of  $U^*$  is

$$V_p = \frac{\sqrt{n}}{(n+1)} \sum_{k=1}^n (x_k - \bar{x}) [F(y_k) - 1/2].$$

Since  $F(y_k)$  has a uniform distribution on  $(0, 1)$  with mean  $1/2$  and variance  $1/12$ , we have

$$\text{Var } V_p = \frac{n}{12(n+1)^2} \sum_{k=1}^n (x_k - \bar{x})^2 \rightarrow \frac{1}{12} \delta^2.$$

Hence  $\text{Var } U^* \rightarrow \delta^2/12$ , so  $E(U^* - V_p)^2 \rightarrow 0$ .

Thus,  $U^*$  and  $V_p$  have the same distribution.

That is,  $V_p$  (and hence  $U^*$ ) is asymptotically  $N(0, \delta^2/12)$ .

### 3. Confidence intervals for $\beta$

3.1. Confidence interval based on the Theil test.

For a symmetric two-sided confidence interval for  $\beta$  with confidence coefficient  $1 - \alpha$ ,

a) Determine the constant  $C_\alpha$  that satisfies the equation

$$P_0 \{-C_\alpha \leq C \leq C_\alpha\} = 1 - \alpha.$$

Note that  $C_\alpha + 2 = k((\alpha/2), n)$ .

This values were evaluated from the Kendall's K-statistic.

b) Obtain the ordered values  $S^{(1)} < \dots < S^{(N)}$  of the  $N = \binom{n}{2}$  sample slopes  $S_{ij} = (Y_j - Y_i) / (x_j - x_i)$ .

c) Set  $M_1 = \frac{N - C_\alpha}{2}$  and  $M_2 = \frac{N + C_\alpha}{2}$ .

d) The  $1 - \alpha$  confidence interval  $(\beta_L, \beta_U)$  is defined by  $\beta_L = S^{(M_1)}$  and  $\beta_U = S^{(M_2)}$ .

Hence we have  $P_{\beta} \{ \beta_L < \beta < \beta_U \} = 1 - \alpha$ .

### 3.2. Confidence interval based on $T_\beta$

Consider the simple linear regression model (1.1) and the test of  $H_0: \beta = 0$  against  $H_1: \beta > 0$ , where the test is based on the statistic  $T_{\beta 0}$ . From the theorem 2, the test which rejects  $H_0$  if  $T_{\beta 0} > (a \dots / 2) + z_\alpha (\sum_{i=1}^n A_i^2 / 12)^{1/2}$  is an approximate level  $\alpha$  test, where  $z_\alpha$  is the quantile of order  $1 - \alpha$  for the standard normal distribution.

A confidence interval for  $\beta$  can be obtained by inverting the two-sided alternatives.

Let  $H(s)$  denote the cumulative distribution function of the discrete probability which assigns probability  $a_i / a \dots$  to the points  $S_{ij}$ . Let  $H^{-1}(u) = \inf \{ s : H(s) \geq u \}$ . Suppose that  $t_1$  and  $t_2$  satisfy  $P_{\beta} (t_1 \leq T_{\beta} < t_2) = 1 - \alpha$ . We note that  $P(T_{\beta} < a \dots / 2) \leq P(\hat{\beta} < t) \leq P(T_{\beta} < a \dots / 2)$ . Since the event  $\{ t_1 \leq T_{\beta} < t_2 \}$  is equivalent to the event  $\{ H^{-1}(t_1) \leq \beta < H^{-1}(t_2) \}$ , the interval  $(H^{-1}(t_1), H^{-1}(t_2))$  is a  $(1 - \alpha)$  100 percent confidence interval for  $\beta$ . By theorem 2, we

determine the constants  $t_1$  and  $t_2$ . An approximate  $(1 - \alpha)$  100 percent confidence interval for  $\beta$ .

$$(H^{-1}((1/2) - Z\alpha/2) (\sum_{i=1}^n A_i^2 / 12)^{1/2} / a \dots).$$

$$(H^{-1}((1/2) + Z\alpha/2) (\sum_{i=1}^n A_i^2 / 12)^{1/2} / a \dots).$$

## 4. Conclusion

A natural estimate of  $\beta$  based on Theil statistic was the median of the pairwise slopes. Now various weights  $a_{ij}$  are introduced.

a) Let the weights be given by  $a_{ij} = 1$ ,  $i < j$ , if  $x_i \neq x_j$ ; otherwise let  $a_{ij} = 0$ .

If the  $x_i$ 's are all distinct, then  $A_i = 2i - (n + 1)$ ,  $a \dots = n(n - 1) / 2$ , and conditions A and B hold.

b) Let  $a_{ij} = j - i$ ,  $i < j$ , if  $x_i \neq x_j$ ; otherwise 0.

c) Let  $a_{ij} = x_j - x_i$ ,  $i < j$ . Then  $A_i = n(x_i - \bar{x})$  and  $T_\beta$  is asymptotically normal with mean  $\sum_{i=1}^n \sum_{j=1}^n (x_j - x_i) / 2$  and variance  $n^2 \sum_{i=1}^n (x_i - \bar{x})^2 / 12$ .

Another weight  $a_{ij}$  can be given by  $a_{ij} = (x_j - x_i) / (j - i)$ ,  $i < j$ , but this weight  $a_{ij}$  does not have a simple representation in evaluating  $T$ .

If the above weights are considered in view of the efficiency considerations, the weight  $a_{ij} = x_j - x_i$  is recommended.

That is, the weight  $a_{ij} = x_j - x_i$  merits serious consideration as an alternative to the classical estimate.

## Literature cited

Adichie, J. N., 1967. Estimates of regression parameters based on rank test. *A. M. S.* 38: 894-904.

Jarosiav Hájek and Zbynek Šidák, 1967. *Theory of rank test.* Academic press, New York pp. 56-165.

Lehmann, E. L., 1975. *Nonparametric statistical method based on ranks.* McGraw-Hill, New York. pp. 82-93.

Myles Hollander and D. A. Wolfe, 1973. *Nonparametric statistical methods.* John Willey & SONS pp. 67-75.

Pranab Kumar Sen, 1981. *Sequential nonparametrics.* John Willey & SONS. pp. 86-90.

Randles R. H. and Wolfe D. A. 1979. *Introduction to the theory of nonparametric statistics.* John Willey & SONS. pp. 103-113.

## 국 문 초 록

순위 검정 통계량의 형태인 무제있는 U-통계량을 사용해서 단순 회귀 모형에 대한 Theil의 검정을 살펴 보며 회귀계수  $\beta$ 에 대한 신뢰구간을 추정한다.