

## The mathematics of Stute (1997) test

Let  $Y$  and  $D$  be two random variables. Let  $m(D) = E[Y|D]$ . The null hypothesis of the test is that  $m(D) = \alpha_0 + \alpha_1 D$  for two real numbers  $\alpha_0$  and  $\alpha_1$ . This means that, under the null,  $m(\cdot)$  is linear in  $D$ . This hypothesis can be tested in a sample with  $N$  i.i.d. realizations of  $(Y, D)$  using the following test statistic from Stute (1997):

$$S = \frac{1}{N^2} \sum_{i=1}^N \left( \sum_{j=1}^i \varepsilon_{(j)} \right)^2$$

where  $\varepsilon_{(j)}$  is the residual from a linear regression of  $Y$  on  $D$  and a constant of the  $j$ -th observation after sorting by  $D$ . In other words,  $S$  is obtained by sorting the data from the smallest to the largest value of  $D$  and summing the squares of the total cumulative sums of the linear regression residuals.

Stute et al. (1998) show that, under the null,  $S$  is finite. Conversely, under the alternative, at least one of the inner sums tends to infinity, hence  $S$  diverges. Inference is performed using wild bootstrap. Specifically,  $S$  is re-computed replacing  $Y$  with  $Y^*$ , i.e. the predicted value of  $Y$  from the linear regression of  $Y$  on  $D$  and a constant, plus the residuals multiplied by a two-point random variable, denoted as  $V_{(j)}$ , such that:

$$P\left(V_{(j)} = \frac{1 + \sqrt{5}}{2}\right) = \frac{\sqrt{5} - 1}{2}, P\left(V_{(j)} = \frac{1 - \sqrt{5}}{2}\right) = \frac{3 - \sqrt{5}}{2}.$$

Denote with  $S_b^*$  the  $S^*$  statistic computed at the  $b$ -th bootstrap replication. The p-value from  $B$  bootstrap replications is computed as

$$\frac{1}{B} \sum_{b=1}^B 1\{S < S_b^*\}$$

Intuitively, under the alternative, the p-value should be zero, due to the fact that  $S$  diverges.

This test also works with panel data. In that case, the  $S$  statistic is computed for each value of the time variable. Moreover,  $V_{(j)}$  remains constant at the group level across the computation of the period-specific test statistics. Hence, the residual of group  $g$  from a linear regression of  $Y_{g,t}$  on  $D_{g,t}$  and a constant are multiplied by  $V_g$ , regardless of  $t$ . Lastly, the individual test results can be

summed into a joint test statistic. In this case, inference is performed using the distribution of the sum of the bootstrap statistics. Denote with  $S_\ell$  the period- $\ell$  test statistic and with  $S_{\ell,b}^*$  its  $b$ -th bootstrap estimate. In a dataset with  $L$  periods, the p-value of the joint test is computed as follows:

$$\frac{1}{B} \sum_{b=1}^B 1 \left\{ \sum_{\ell=1}^L S_\ell < \sum_{\ell=1}^L S_{\ell,b}^* \right\}.$$

## References

- Stute, W. (1997). Nonparametric model checks for regression. *The Annals of Statistics*.
- Stute, W., W. G. Manteiga, and M. P. Quindimil (1998). Bootstrap approximations in model checks for regression. *Journal of the American Statistical Association*.