# CERTAIN DISTRIBUTION-FREE TESTS OF REGRESSION*

By R. S. Kurup ‡

*University of Kerala, Trivandrum*

## 1. The Problem

SUPPOSE we are given $n$ pairs of observations $(x_i, y_i)$, $i = 1, 2 \cdots n$ from a continuous bivariate distribution and we are required to fit a relation of the form $Y = f(x, \theta)$ where '$\theta$' denotes a set of parameters whose values may be found by any method of estimation. To test the significance of regression, the null hypothesis is $H_0 : \theta = 0$. Classical workers tested regression by assuming that the errors are normally and independently distributed and this forms the basis of the $x^2$-test. In this paper the problem is tackled without any such assumptions.

For this problem, Brown and Mood (1950),[1] (1951)[2] suggested a statistic,

$$A = \frac{8}{n} \left\{ \left( r_1 - \frac{n}{4} \right)^2 + \left( r_2 - \frac{n}{4} \right)^2 \right\}$$

where $r_1$ and $r_2$ are the number of positive $\epsilon$'s below and above the median of the $x$'s, $\epsilon$ being the discrepancy between the observed '$y$' and the value of '$y$' under the null hypothesis. For moderately large '$n$', this is distributed as a '$x^2$' with 2 degrees of freedom. This statistic considers the 4 possible arrangements of signs, as shown below:

| $n/2$ | | | | | | $n/2$ | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| + | + | + | + | $\cdots$ | + | − | − | − | − | $\cdots$ | − |
| + | + | + | + | $\cdots$ | + | + | + | + | + | $\cdots$ | + |
| − | − | − | − | $\cdots$ | − | − | − | − | − | $\cdots$ | − |
| − | − | − | − | $\cdots$ | − | + | + | + | + | $\cdots$ | + |

Daniels in 1954[3] suggested:

$$B = \frac{2}{\sqrt{n}} \left\{ \left| r_1 - \frac{n}{4} \right| + \left| r_2 - \frac{n}{4} \right| \right\}$$

as a test criterion. This has the asymptotic distribution:

$$P(B \geqslant B_0) = 4 \phi(B_0) \left(1 - \phi(B_0)\right)$$

where $\phi(B_0)$ is the normal cumulative distribution. It was shown that this is more powerful than the $A$-test mentioned earlier. Daniels proposed another test, the $m$-test, based on the $2^n$ possible arrangements of signs.

## 2. THE PROPOSED TEST-CRITERIA AND THEIR DISTRIBUTIONS

Let $x_1, x_2 \cdots x_n$ denote the ordered $x$'s in ascending order of magnitude; also let $\epsilon$ be the difference between the observed '$y$' and the value of '$y$' under the null hypothesis and '$R_i$' denote the number of positive $\epsilon$'s up to $x_i$ (including '$x_i$'). Considering only four possible arrangements of the signs we may formulate the criterion $R_n$ with expectation $n/2$ and variance $n/4$. It is obvious that for large $n$, $4/n(R_n - n/2)$ follows the normal distribution with mean '0' and variance unity under '$H_0$'. For small '$n$', the probability of any particular '$R_n$' can be found easily by computation as the probability is $\frac{1}{2}$ for any $\epsilon$ to be positive or negative.

Another criterion which can be used for the test is $T = \sum\limits_{i=1}^{n} R_i$.

The exact distribution of this criterion has been tabulated up to '$n = 10$' by considering all possible arrangements. As '$T$' is found to be symmetrical about its mean $n(n+1)/4$ the upper half alone is given in Table I.

The significance of '$T$' should be tested using the two tails of the distribution.

The distribution of $T$ is symmetrical. The mean value of '$T$' is $n(n+1)/4$ and variance $[n(n+1)(2n+1)]/24$

$$T' = \frac{T - \dfrac{n(n+1)}{4}}{\sqrt{\dfrac{n(n+1)(2n+1)}{24}}}$$

TABLE I

*Ordinates $P$ of the distribution of $T = \sum\limits_{i=1}^{n} R_i$*

| $n$ | $T$ | $P$ | $n$ | $T$ | $P$ |
|---|---|---|---|---|---|
| 3 | 6 | ·125 | | 22 | ·03125 |
| | 5 | ·125 | | 21 | ·0390625 |
| | 4 | ·125 | | 20 | ,, |
| | 3 | ·250 | | 19 | ·046875 |
| | | | | 18 | ·0546875 |
| 4 | 10 | ·0625 | | 17 | ,, |
| | 9 | ,, | | 16 | ·0625 |
| | 8 | ,, | | 15 | ,, |
| | 7 | ·1250 | | 14 | ,, |
| | 6 | ,, | | | |
| | 5 | ,, | 8 | 36 | ·00390625 |
| | | | | 35 | ,, |
| 5 | 15 | ·03125 | | 34 | ,, |
| | 14 | ,, | | 33 | ·0078125 |
| | 13 | ,, | | 32 | ,, |
| | 12 | ·0625 | | 31 | ·01171875 |
| | 11 | ,, | | 30 | ·015625 |
| | 10 | ·09375 | | 29 | ·01953125 |
| | 9 | ,, | | 28 | ·0234375 |
| | 8 | ,, | | 27 | ·02734375 |
| | | | | 26 | ·03125 |
| 6 | 21 | ·015625 | | 25 | ·03515625 |
| | 20 | ,, | | 24 | .0390625 |
| | 19 | ,, | | 23 | ·04296875 |
| | 18 | ·03125 | | 22 | ·046875 |
| | 17 | ,, | | 21 | ·05078125 |
| | 16 | ·046875 | | 20 | ,, |
| | 15 | ·0625 | | 19 | ,, |
| | 14 | ,, | | 18 | ·0546875 |
| | 13 | ,, | | | |
| | 12 | ·078125 | 9 | 45 | ·001953125 |
| | 11 | ,, | | 44 | ,, |
| | | | | 43 | ,, |
| 7 | 28 | ·0078125 | | 42 | ·00390625 |
| | 27 | ,, | | 41 | ,, |
| | 26 | ,, | | 40 | ·005859375 |
| | 25 | ·015625 | | 39 | ·0078125 |
| | 24 | ,, | | 38 | ·009765625 |
| | 23 | ·0234375 | | 37 | ·01171875 |

TABLE I (*Contd.*)

| $n$ | $T$ | $P$ | $n$ | $T$ | $P$ |
|---|---|---|---|---|---|
|  | 36 | ·015625 |  | 49 | ·00390625 |
|  | 35 | ·017578125 |  | 48 | ·0048828125 |
|  | 34 | ·01953125 |  | 47 | ·005859375 |
|  | 33 | ·0234375 |  | 46 | ·0078125 |
|  | 32 | ·025390625 |  | 45 | ·009765625 |
|  | 31 | ·029296875 |  | 44 | ·0107421875 |
|  | 30 | ·033203125 |  | 43 | ·0126953125 |
|  | 29 | ·03515625 |  | 42 | ·0146484375 |
|  | 28 | ·037109375 |  | 41 | ·0166015625 |
|  | 27 | ·041015625 |  | 40 | ·01953125 |
|  | 26 | ,, |  | 39 | ·021484375 |
|  | 25 | ·04296875 |  | 38 | ·0234375 |
|  | 24 | ·044921875 |  | 37 | ·0263671875 |
|  | 23 | ,, |  | 36 | ·0283203125 |
|  |  |  |  | 35 | ·0302734375 |
| 10 | 55 | ·0009765625 |  | 34 | ·0322265625 |
|  | 54 | ,, |  | 33 | ·0341796875 |
|  | 53 | ,, |  | 32 | ·03515625 |
|  | 52 | ·001953125 |  | 31 | ·037109375 |
|  | 51 | ,, |  | 30 | ·0380859375 |
|  | 50 | ·0029296875 |  | 29 | ,, |
|  |  |  |  | 28 | ·0390625 |

(,,) Denotes the same value as before.

will therefore tend to have a normal distribution with mean '0' and variance 1 for large '$n$'.

### 3. ANOTHER TEST CRITERION

If
$$z_i = R_i - \frac{i}{2},$$

$$E(z_i) = 0$$

and

$$E(z_i z_j) = \frac{n_j}{4} \quad \text{for} \quad j \leqslant i$$

$$= \frac{n_i}{4} \quad \text{for } j > i.$$

Denoting by $z$, the matrix of values $z_1, z_2, \cdots z_n$ and by $\Gamma$ the covariance matrix of the $z_i$'s

$$T_n^2 = z\,\Gamma^{-1}z'$$

$$= \frac{4}{n}\left[2\sum_{i=1}^{n-1}(z_i^2 - z_iz_{i+1}) + z_n^2\right]$$

$$= \frac{4}{n}\sum_{i=1}^{n}\left(\delta_i - \frac{1}{2}\right)^2$$

where $\delta_i$ is 1 or 0 with probability $\frac{1}{2}$. Here, it may be noted that $R_i$ (defined earlier) $= \sum_{j=1}^{i}\delta_j$.

From the structure of the criterion $T_n^2 = z\Gamma^{-1}z'$, it can be seen that for large $n$, it will behave as a $x^2$ with '$n$' degrees of freedom. For small '$n$' the exact distribution under $H_0$ can be tabulated.

## 4. POWER OF THE CRITERIA

We may for convenience assume that we are required to test the linear relation $y = a + \beta x$. Under the alternative $(\beta, a)$ against $(\beta_0, a_0)$, $\epsilon_i = y_i - a - \beta x_i \cdots$ (1)

has probability $\frac{1}{2}$ for being positive or negative

$$y_i - a_0 - \beta_0 x_i = \epsilon_i - (a_0 - a) - (\beta_0 - \beta)\,x_i. \tag{2}$$

Let

$$p_i = \text{prob. }\{\epsilon_i > (a_0 - a) + (\beta_0 - \beta)\,x_i\}. \tag{3}$$

Following Daniels[3] we may consider $p_i = p$ for all '$i$'. This is the case if the alternatives are $(\beta_0, a)$ and all the $\epsilon$'s have the same distribution $f(\epsilon)$. Then

$$p = \int_{a_0-a}^{a} f(\epsilon)\,d\epsilon \simeq \frac{1}{2} - (a_0 - a)f(0) \tag{4}$$

If as in [3]

$$p = \frac{1}{2}\left(1 - \frac{\mu}{\sqrt{n}}\right),$$

we have

$$\mu = 2(a_0 - a)f(0)\sqrt{n} \tag{5}$$

The limiting form of the distribution of $R_n$ under the alternative hypothesis is a normal distribution with parameter

$$\frac{np - \frac{n}{2}}{\sqrt{npq}} \simeq \mu$$

for large '$n$'. Using this, the power of $R_n$ for large '$n$' has been computed and given in Table II for various values of '$\mu$'. For comparison the relevant results from [3] are reproduced.

TABLE II

*Asymptotic power of five tests at ·05 level for alternatives $a \neq a_0$, $\beta = \beta_0$*
*Here $\mu = \sqrt{n}\{1 - 2 \text{ prob. } (\epsilon_i > (a - a_0))\}$*

| $\mu$ | $m$-Test | $A$ | $B$ | $R_n$ | $F$-Test |
|-------|----------|-----|-----|-------|----------|
| ·0 | ·05 | ·05 | ·05 | ·05 | ·05 |
| ·790 | ·09 | ·10 | ·10 | ·12 | ·13 |
| 1·316 | ·16 | ·20 | ·20 | ·26 | ·29 |
| 1·666 | ·25 | ·30 | ·30 | ·39 | ·45 |
| 1·958 | ·33 | ·40 | ·41 | ·50 | ·59 |
| 2·226 | ·42 | ·50 | ·51 | ·61 | ·71 |
| 2·493 | ·52 | ·60 | ·61 | ·70 | ·81 |
| 2·775 | ·62 | ·70 | ·71 | ·79 | ·89 |
| 3·104 | ·73 | ·80 | ·81 | ·87 | ·95 |
| 3·557 | ·85 | ·90 | ·91 | ·95 | ·99 |

The values given above show that $R_n$ is more powerful than the other distribution-free tests namely $m$-test, $A$ and $B$. Also it approaches the $F$-Test (in power) based on the assumption of normality (parameter $\frac{1}{2}\pi\mu^2$ and degrees of freedom 2).

Under the alternative hypothesis discussed above the criterion $T$ has expectation $[n(n + 1)] P/2$ and variance $[n(n + 1)] pq/2$. Hence in terms of $\mu$ the parameter for $T$ is

$$r = -\sqrt{\frac{n(n+1)\mu^2}{2(n-\mu^2)}}$$

and $T$ is asymptotically normally distributed. For the criterion $T_n^2$ the corresponding parameter is $r^2$. The power of $T$ and $T_n^2$ is not evaluated here as they behave like the '$t$' and $x^2$-tests respectively.

## REFERENCES

1. Mood, A. M.    ..    *Introduction to the Theory of Statistics*, McGraw-Hill, New York, 1950.

2. Brown, G. W. and Mood, A. M.    "On median tests of linear hypothesis," *Proceedings of the Second Berkely Symposium on Mathematical Statistics and Probability*, University of California, Berkeley, 1951, pp. 159–66.

3. Daniels, H. E.    ..    "A distribution-free test for regression parameters," *Annals of Math-Statistics*, 1954, **25**, 499–514.