

# $\chi^2$ test

## Testing the goodness of fit of a geometric model

In S2 Section 4.2 you met random numbers as a device for selecting random samples. Below is the start of a list of digits obtained by using the random number generator on a calculator. If such a list is random, then each digit is equally likely to appear in any position in the list. As a result odd and even digits are also equally likely to appear in any position on the list.

6 3 3 0 4 5 3 2 1 4 1 3 2 1 2 7 2 2 5 1 1 1 4 1 2 5 3...

One way of testing for the property of randomness is to look at the intervals separating the odd digits. The list is repeated below with a line placed after each odd digit. This has the effect of breaking the sequence up into 'runs', consisting of even digits, until the last digit, which is odd. The lengths of the runs below are 2, 1, 3, 1, 2, 2, etc.

6 3|3|0 4 5|3|2 1|4 1|3|2 1|2 7|2 2 5|1|1|1|4 1|2 5|3|...

If the digits are random, then the length,  $Y$ , of a run will have a geometric distribution. You can see this by considering an even number as a 'failure' and an odd number as a 'success'. In this case  $p = q = \frac{1}{2}$ , so  $Y \sim \text{Geo}(\frac{1}{2})$ . Table 5.11 gives the frequency distribution of run length for 50 runs.

Length of run	1	2	3	4	5	6	7	>7
Frequency	21	18	5	2	3	0	1	0

Table 5.11. Frequency distribution of 'runs' of even digits.

A  $\chi^2$  goodness of fit test can be used to test whether the data in Table 5.11 can be modelled by the distribution  $\text{Geo}(\frac{1}{2})$ . The null and alternative hypotheses are

$H_0$ : the data can be modelled by  $\text{Geo}(\frac{1}{2})$ ;

$H_1$ : the data cannot be modelled by  $\text{Geo}(\frac{1}{2})$ .

The geometric probabilities are given by  $P(Y = y) = (\frac{1}{2})^{y-1} (\frac{1}{2}) = (\frac{1}{2})^y$  and the expected frequencies are found by multiplying these probabilities by the total observed frequency of 50. For example, the expected frequency for a run length of 2 is equal to  $(\frac{1}{2})^2 \times 50 = 12.5$ .

Unlike the binomial distribution there is no upper limit on the value of  $Y$ . Although the longest run length in Table 5.11 is 7, greater run lengths are possible. For this reason a table of the expected and observed frequencies must also include a class for values greater than 7. The expected frequency for this class must be such that the total expected frequency is 100 and so its value (in this case 0.39) can be found by subtraction. Table 5.12 shows the calculation of  $X^2$ . Classes have been combined where necessary to make all expected frequencies at least 5.

Run length	$O_i$	$E_i$	$O_i - E_i$	$(O_i - E_i)^2 / E_i$
1	21	25	-4	0.64
2	18	12.5	+5.5	2.42
3	5	6.25	-1.25	0.25
4	2	3.13		
5	3	1.56		
6	0	0.78		
7	1	0.39		
>7	0	0.39		
Total	50	50	0	$X^2 = 3.32$

Table 5.12. Calculation of  $X^2$  for the data in Table 5.11.

There are 4 classes after combination and 1 constraint as the total observed frequency of 50 has been used to calculate the expected frequencies. So  $\nu = 4 - 1 = 3$ . From the table on page 303, the rejection region for a 5% significance level is  $X^2 > 7.815$ . Since 3.32 is not in the rejection region, the null hypothesis that the model  $\text{Geo}(\frac{1}{2})$  is suitable is accepted.

- 12 A survey is carried out at a supermarket till. When the till opens, the number of customers up to and including the first person to use one of the carrier bags provided by the supermarket is recorded. This is repeated on 100 consecutive days. The data are summarised in the table below.

Number of customers	1	2	3	4	>4
Frequency	79	15	3	3	0

It is thought that this distribution may be modelled by a geometric distribution with parameter  $p$ , where  $p$  is the probability that a person uses a supermarket carrier bag.

- (a) Calculate the mean and hence obtain an estimate of  $p$ .  
(b) Carry out a test, at the 5% significance level, of the goodness of fit of the model to the data.

Ans

a) 1.3, 0.769

b)  $\text{Geo}(0.769)$  ; 76.92, 17.75, 4.10,  
0.95, 0.28 ; 0.57 ; accept  $H_0$ .