

Goodness-of-fit tests for fitted RRs

Author : P.N. Lee

Date : 11th June 2012

1. Prospective studies

We have an observed table of pseudo-numbers:

<u>Level</u>	<u>Cases</u>	<u>At risk</u>
Baseline	$A_0$	$N_0$
Low exposure	$A_1$	$N_1$
High exposure	$A_2$	$N_2$
Total	$A_S$	$N_S$

We have fitted a set of RRs,  $1 : R_1 : R_2$  ( $R_0 = 1$ )

We wish to derive a set of fitted cases  $F_0, F_1, F_2$

We have the following formulae:

$$F_0 + F_1 + F_2 = A_S \quad (\text{marginal totals stay the same}) \quad (1)$$

$$R_1 = (F_1 N_0) / (F_0 N_1) \quad (2)$$

$$R_2 = (F_2 N_0) / (F_0 N_2) \quad (3)$$

$$\text{From (2)} \quad F_1 = F_0 N_1 R_1 / N_0 \quad (4)$$

$$\text{From (3)} \quad F_2 = F_0 N_2 R_2 / N_0 \quad (5)$$

$$\text{From (1,4,5)} \quad F_0 + \frac{F_0 N_1 R_1}{N_0} + \frac{F_0 N_2 R_2}{N_0} = A_s \quad (6)$$

$$\text{so} \quad F_0 N_0 + F_0 N_1 R_1 + F_0 N_2 R_2 = A_s N_0 \quad (7)$$

$$\text{or} \quad F_0 = (A_s N_0 R_0) / \sum_{i=0}^2 (N_i R_i) \quad (8)$$

$$\text{From (4,8)} \quad F_1 = (A_s N_1 R_1) / \sum_{i=0}^2 (N_i R_i) \quad (9)$$

$$\text{From (5,8)} \quad F_2 = (A_s N_2 R_2) / \sum_{i=0}^2 (N_i R_i) \quad (10)$$

This allows one to derive fitted values and is clearly generalisable to multiple exposure levels (k).

A chisquared test of goodness-of-fit may be derived in the usual way from the formula:

$$\chi^2 = \sum_{i=0}^k (A_i - F_i)^2 / F_i$$

The degrees of freedom is k – 1.

## 2. Case-control studies

Here the observed table of pseudo-numbers is:

<u>Level</u>	<u>Cases</u>	<u>Controls</u>	<u>Total</u>
Baseline	A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>
Level 1	A <sub>1</sub>	B <sub>1</sub>	C <sub>1</sub>
Level 2	A <sub>2</sub>	B <sub>2</sub>	C <sub>2</sub>
Total	A <sub>s</sub>	B <sub>s</sub>	C <sub>s</sub>

The expected table of fitted numbers is:

<u>Level</u>	<u>Cases</u>	<u>Controls</u>
Baseline	$F_0$	$G_0$
Level 1	$F_1$	$G_1$
Level 2	$F_2$	$G_2$

We have fitted RRs,  $1 : R_1 : R_2$  ( $R_0 = 1$ )

### 2.1 A first attempt to solve the equations

The reader may prefer to skip to section 2.2.

We can write down the following formulae based on the marginal totals and the RRs:

$$F_0 + G_0 = C_0 \quad (1)$$

$$F_1 + G_1 = C_1 \quad (2)$$

$$F_2 + G_2 = C_2 \quad (3)$$

$$F_0 + F_1 + F_2 = A_s \quad (4)$$

$$R_1 = F_1 G_0 / (F_0 G_1) \quad (5)$$

$$R_2 = F_2 G_0 / (F_0 G_2) \quad (6)$$

$$\text{From (5)} \quad G_1 = F_1 G_0 / (F_0 R_1) \quad (7)$$

$$\text{From (6)} \quad G_2 = F_2 G_0 / (F_0 R_2) \quad (8)$$

$$\text{From (2,5)} \quad F_1 + \frac{F_1 G_0}{F_0 R_1} = C_1 \quad (9)$$

$$\text{or} \quad F_0 F_1 R_1 + F_1 G_0 = C_1 F_0 R_1 \quad (10)$$

$$\text{From (1)} \quad F_0 F_1 R_1 + F_1 (C_0 - F_0) = C_1 F_0 R_1 \quad (11)$$

$$\text{or} \quad F_1 = C_1 F_0 R_1 / (F_0 R_1 + C_0 - F_0) \quad (12)$$

$$\text{Similarly} \quad F_2 = C_2 F_0 R_2 / (F_0 R_2 + C_0 - F_0) \quad (13)$$

$$\text{From (4,12,13)} \quad F_0 + \frac{C_1 F_0 R_1}{(F_0 R_1 + C_0 - F_0)} + \frac{C_2 F_0 R_2}{(F_0 R_2 + C_0 - F_0)} = A_s \quad (14)$$

This is an equation in  $F_0$  only, which hopefully can be solved.

Formula (12) gives  $F_1$  in terms of  $F_0$

Formula (13) gives  $F_2$  in terms of  $F_0$

Formulae (1,2,3) then give  $G_i$  in terms of  $F_i$

This gives the whole table of fitted numbers. Here we can calculate a goodness-of-fit statistic as:

$$\chi^2 = \sum_{i=0}^k (A_i - F_i)^2 / F_i + \sum_{i=0}^k (B_i - G_i)^2 / G_i$$

The degrees of freedom is  $2k - 1$ .

The problem is solving formula (14). For  $k = 2$  it is a cubic (as can be seen by multiplying through by the denominators). For general  $k$  it involves powers of  $k + 1$  presumably, and I don't know how to solve that. I suppose one could search starting at  $F_0 = A_0$ .

## 2.2 An alternative approach

This is illustrated in T:/PNLEE/FITFORCCSTUDIES.XLSX. Here we have an observed table of pseudo-numbers of

<u>Level</u>	<u>Cases</u>	<u>Controls</u>	<u>Total</u>
Baseline	$A_0 = 40$	$B_0 = 110$	$C_0 = 150$
Level 1	$A_1 = 90$	$B_1 = 110$	$C_1 = 200$
Level 2	$A_2 = 220$	$B_2 = 80$	$C_2 = 300$
Total	$A_S = 350$	$B_S = 300$	$C_S = 650$

We have fitted RRs at 1 : 2 : 4

In line 7 of the spreadsheet we first calculate  $B_i R_i$  for each level and in total, and then derive a first estimate of fitted cases using the formula

$$F_i = (A_S B_i R_i) / \sum_{i=0}^e (B_i R_i)$$

Taking the first estimate of fitted controls ( $G_i$ ) as the original pseudo number ( $B_i$ ), we then have a set of cases and controls which produce the fitted RRs but with the marginal totals not adding up over level.

In line 8, we scale the estimates of  $F_i$  and  $G_i$  by the same factor so that their sum is equal to  $C_i$ . Now the marginal totals over all cases and over all controls do not equal  $A_S$  and  $B_S$ . So in line 9, we scale them so that they do sum to  $A_S$  and  $B_S$ . However, the totals over levels,  $C_i$ , are then wrong.

We therefore repeat the two steps until they converge to the right answer, which is

<u>Level</u>	<u>Cases</u>	<u>Controls</u>	<u>Total</u>
Baseline	$F_0 = 50$	$G_0 = 100$	$C_0 = 150$
Level 1	$F_1 = 100$	$G_1 = 100$	$C_1 = 200$
Level 2	$F_2 = 200$	$G_2 = 100$	$C_2 = 300$
Total	$A_5 = 350$	$B_5 = 300$	$C_5 = 650$

This happens rapidly as can be seen and by 5 goes (line 16) the fit is very good.

Lines 22-35 in the spreadsheet correspond to lines 7 – 20 and show the chisquared statistic.