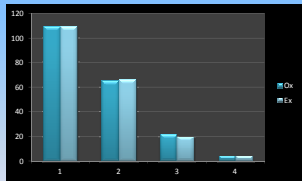


# STATS 101 Introductory Statistics CHI SQUARE TESTS



www.stats24x7.com

1

• **CHI-SQUARE TESTS** are used for

1. Testing goodness of fit
2. Testing independence of two ATTRIBUTES in a table
3. Testing homogeneity of proportions in a table

In this lecture, you will learn how to run chi-square tests in EXCEL (see chi square tests.xlsx)

www.stats24x7.com

2

## The Chi-Square Test

In each application of the chi-square test, we will be given –

$O_j$  = observed frequency of cell j

$E_j$  = frequency of cell j computed assuming null hypothesis to be true

The idea behind the chi-square test is simple – if all of the expected frequencies  $E_j$  are 'close' to the corresponding observed frequencies  $O_j$  then the null hypothesis should be true. The chi-square test, therefore, has the following form:

www.stats24x7.com

3

Reject  $H_0$  if the *DISTANCE* between the observed and expected frequencies is *LARGE*. The chi-square distance between the observed and expected frequencies ( $O$  and  $E$ ) given below is used:

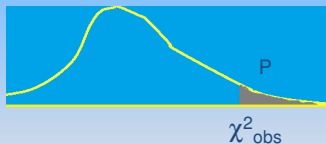
$$DIST(O - E) = \chi_{obs}^2 = \sum \frac{(O_j - E_j)^2}{E_j}$$

www.stats24x7.com

4

In each application of the chi-square test, the null is rejected if the P-value:

$$P = P(\chi_{df}^2 > \chi_{obs}^2) < .05 \text{ if using a test of size } .05$$



The degree of freedom (df) of the chi-square test depends on the chi-square problem.

www.stats24x7.com

5

## 1. Testing Goodness of Fit (GOF) of a specified probability

• If data is continuous, calculate a histogram (this step is not needed for discrete data); this gives us OBSERVED FREQUENCIES  $O_i$  for each class interval.

• Calculate EXPECTED FREQUENCY  $E_i$  for Cell (class interval)  $i$  by the formula

$$E_i = N \times P(\text{an observation falls in Cell } i \text{ if } H_0 \text{ is true}).$$

• Calculate  $\chi_{obs}^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$

Here  $k = \#$  of cells or class intervals

www.stats24x7.com

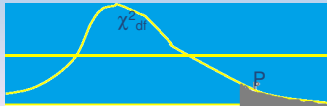
6

Reject  $H_0$  (conclude data does not follow hypothesized probability distribution) if

$$\chi^2_{obs} > \chi^2_{df}$$

$$df = k - \# \text{ of parameters estimated} - 1$$

Alternatively, reject  $H_0$  if the P-value as shown in following figure is  $< .05$



www.stats24x7.com

7

### TESTING GOODNESS OF FIT EXAMPLE (Discrete data)

Ex 1: Use the data (x, O) to test if the die is fair:

x	O <sub>i</sub>	E <sub>i</sub>	chi-sqr
1	28	30	0.13
2	36	30	1.2
3	36	30	1.2
4	30	30	0
5	27	30	0.3
6	23	30	1.63
SUM	180	180	4.46

In this table:

O<sub>x</sub> = observed frequency of value x  
E<sub>x</sub> = EXPECTED frequency of value x

E<sub>x</sub> = frequency expected under the null hypothesis  $H_0$

$H_0: f(x) = 1/6, x = 1, 2, \dots, 6$

$E_x = N \times f(x) = 180 (1/6) = 30$  for each x

$\chi^2$  column =  $(O_x - E_x)^2 / E_x$

www.stats24x7.com

8

$$H_0: P(X = x) = 1/6, x = 1, 2, \dots, 6$$

$$H_1: H_0 \text{ - false}$$

FACT: If the null hypothesis is true,

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{df}$$

$$df = k - 1 - \#(\text{parameters estimated})$$

$df = k - 1 = 6 - 1 = 5$  for the above example, as the null hypothesis specified the probabilities and no parameters had to be estimated (i.e.,  $k = 0$ )

www.stats24x7.com

9

$$\chi^2_{calc} = \frac{(-2)^2 + (6)^2 + (6)^2 + (0)^2 + (-3)^2 + (-7)^2}{30} = \frac{134}{30} = 4.47$$

$$df = 6 - 1 = 5$$

$$\chi^2_{3,95} = 11.07$$

We can get the P-value  $P(\chi^2_5 > 4.47)$  from EXCEL in 2 ways:

1) =chidist(.47,5) will return the P-value = 0.48

2) chitest(observed range, expected range) will only return the P-value, and not  $\chi^2_{obs}$ .

We do not reject the null hypothesis, and conclude that the die is fair.

NOTE: For the above chi-square test to be valid, each  $E_i \geq 5$ .

If this is not the case, then combine 2 or more adjacent cells.

www.stats24x7.com

10

### Using EXCEL to run Chi-square Goodness of Fit (GOF) test for Example 1

- 1) Type observed frequencies in column A, cells 2 – 7. Label the column O<sub>x</sub> (use cell A1).
- 2) In cell C2, type =180/6 to get E<sub>x</sub>.
- 3) Copy formula in cell B2 to cells B3 – B7.
- 4) In cell C1, type 'Chi-sqr' to label the chi-square column.
- 5) In cell C2, type =(A2 – B2)^2/B2 to calculate the chi-square contribution from the 1<sup>st</sup> cell. Copy formula in cell C2 to cells C3 – C7. You now have the table on the next slide.

www.stats24x7.com

11

### Using EXCEL to run GOF test for Example 1

#### EXCEL COLUMN

A	B	C
O <sub>x</sub>	E <sub>x</sub>	Chi-square
28	30	0.13
36	30	1.20
36	30	1.20
30	30	0.00
27	30	0.30
23	30	1.63
TOTAL	180	4.47

To get the P-value for the observed chi-square value, type =CHIDIST(4.47,5) in a blank cell in EXCEL. This will return P = 0.48

Since  $P > .05$ , we conclude that die is fair.

=sum(A2:A7)

www.stats24x7.com

12

## Using EXCEL to run GOF test for Example 1

We can also get the P-value from EXCEL directly, WITHOUT CALCULATING the  $\chi^2_{obs}$  value.

Go to a blank cell in EXCEL, and type

`=chitest(A2:A7,B2:B7)`

This will just yield the P-value of 0.48 (same as on the previous slide).

www.stats24x7.com

13

## 2. CHI-SQUARE TEST OF INDEPENDENCE Example 2: Test if AGE and GAME TYPE are independent, at test size $\alpha = 0.01$ .

GAME	AGE			TOTAL
	21-25	26-50	Over 50	
Multi-Line Slots	15	37	16	68
Video Poker	25	25	17	67
Wheel of Fortune	14	40	27	81
Sports Book	11	4	1	16
Blackjack	9	23	14	46
Megabucks	3	8	1	12
TOTAL	77	137	76	290

www.stats24x7.com

14

Recall that A and B are independent if  $P(A \cap B) = P(A)P(B)$

If  $H_0$  is true:

$$P(\text{Multi line slot, Age 21 - 25}) =$$

$$P(\text{Multi line slot}) \times P(\text{Age 21 - 25})$$

$$= (68/290)(77/290)$$

$$E_{1,1} = 290 \times (68/290)(77/290)$$

$$= 68 \times 77/290$$

Hence

$$E_{i,j} = \frac{\text{Total of Row-i} \times \text{Total of Column j}}{\text{Grand Total}}$$

www.stats24x7.com

15

To run the chi-square test of independence, type the data table and copy row and column headings next to data table as

A	B	C	D	E	F	G	H	I	J	K
AGE				AGE						
GAME	21-25	26-50	Over 50	TOTAL	GAME	21-25	26-50	Over 50	TOTAL	
Multi-Line Slots	15	37	16	68	Multi-Line Slots	<code>=B\$9*\$E3/\$E\$9</code>				
Video Poker	25	25	17	67	Video Poker					
Wheel of Fortune	14	40	27	81	Wheel of Fortune					
Sports Book	11	4	1	16	Sports Book					
Blackjack	9	23	14	46	Blackjack					
Megabucks	3	8	1	12	Megabucks					
TOTAL	77	137	76	290	TOTAL					
OBSERVED FREQUENCIES				EXPECTED FREQUENCIES						

If you typed data exactly as shown, type `=B$9*$E3/$E$9` in cell H3, and then copy this formula in cells I3, J3, ..., H8, I8, J8. This should result in EXPECTED FREQUENCIES table.

www.stats24x7.com

16

## Chi-square test of independence in EXCEL for Example 2 (contd.)

A	B	C	D	E	F	G	H	I	J	K
AGE				AGE						
GAME	21-25	26-50	Over 50	TOTAL	GAME	21-25	26-50	Over 50	TOTAL	
Multi-Line Slots	15	37	16	68	Multi-Line Slots	18.06	32.12	17.82	68	
Video Poker	25	25	17	67	Video Poker	17.79	31.65	17.56	67	
Wheel of Fortune	14	40	27	81	Wheel of Fortune	21.51	38.27	21.23	81	
Sports Book	11	4	1	16	Sports Book	4.25	7.56	4.19	16	
Blackjack	9	23	14	46	Blackjack	12.21	21.73	12.06	46	
Megabucks	3	8	1	12	Megabucks	3.19	5.67	3.14	12	
TOTAL	77	137	76	290	TOTAL	77	137	76	290	
P-value				<code>=CHITEST(B3:D8,H3:J8)</code>	16:D21,I16:K21					
Reject null, conclude Voting Preference and Gender are DEPENDENT or ASSOCIATED.										

In Excel, type `=chitest(B3:D8,H3:J8)`, this will result in the P-value of  $0.0015 < .05$ , and we conclude that GAME PREFERENCE does depend on AGE GROUP.

www.stats24x7.com

17

## 3. CHI-SQUARE TEST OF HOMOGENEITY

Example 3. In a telephone survey, respondents were asked to indicate their level of agreement with the statement "Cigarette smoking should be banned in public places". The results are shown in the table below: SA = strongly agree, A = agree, N = neutral, D = disagree, SD = strongly disagree.

	SA	A	N	D	SD	TOTAL
F	40	38	16	37	5	136
M	16	25	11	25	11	88
TOTAL	56	63	27	62	16	224

Test if there is no difference in Males and Females with respect to their levels of agreement on the banning of smoking in public places. Let  $\alpha = 0.05$ .

www.stats24x7.com

18

$$H_0: p_{Fi} = p_{Mi}, \quad i = 1, 2, 3, 4, 5$$

$H_1$ :  $H_0$  is false

In words: when  $H_0$  is true, there is no difference in Males and Females with respect to their levels of agreement on the banning of smoking in public places.

When  $H_0$  is TRUE:

$$p_{F1} = p_{M1} = P(\text{Strong Agreement, M or F}) = 56/224$$

$$O_{1,1} = \# \text{ of females expected to strongly agree when } H_0 \text{ TRUE} \\ = \text{Total Females} \times P(\text{Strong Agreement, M or F}) = 136 \times 56/224$$

$$O_{2,1} = \# \text{ of Males expected to strongly agree when } H_0 \text{ TRUE} \\ = \text{Total Males} \times P(\text{Strong Agreement, M or F}) = 88 \times 56/224, \\ \text{etc.}$$

www.stats24x7.com

19

Hence

$$O_{i,j} = \frac{\text{Total of Row-}i \times \text{Total of Column } j}{\text{Grand Total}}$$

which is exactly the same formula as in the case of chi-square test for independence. We can therefore run the test of homogeneity in exactly the same we ran the chi-square test of independence.

www.stats24x7.com

20

### EXAMPLE 3 - CHI-SQUARE TEST OF HOMOGENEITY in EXCEL

To run the chi-square test of homogeneity, type the data table and copy row and column headings next to data table as shown here.

A	B	C	D	E	F	G
OBSERVED FREQUENCIES						
Gender	SA	A	N	D	SD	TOTAL
F	40	38	16	37	5	136
M	16	25	11	25	11	88
TOTAL	56	63	27	62	16	224
CIES = ROW TOTAL x COLUMN TOTAL/GRAND TOTAL						
Gender	SA	A	N	D	SD	TOTAL
F						
M						
TOTAL						

www.stats24x7.com

21

In cell B9, type `=B$5*$G3/$G$5`, copy this formula in all other cells to get EXPECTED FREQUENCIES, calculate ROW TOTAL and COLUMN TOTALS in EXCEL. TYPE `=CHITEST(B3:F4,B9:F10)` in cell B14, which gives  $P=.072652$ .

A	B	C	D	E	F	G
OBSERVED FREQUENCIES						
Gender	SA	A	N	D	SD	TOTAL
F	40	38	16	37	5	136
M	16	25	11	25	11	88
TOTAL	56	63	27	62	16	224
CIES = ROW TOTAL x COLUMN TOTAL/GRAND TOTAL						
Gender	SA	A	N	D	SD	TOTAL
F	<code>=B\$5*\$G3/\$G\$5</code>	24.75	10.61	24.36	6.29	88
M	22.00	24.75	10.61	24.36	6.29	88
TOTAL	56	63	27	62	16	224
P	0.072652	> .05				
	DO NOT REJECT NULL					

www.stats24x7.com

22

Since  $P > .05$ , Null is not rejected.  
CONCLUSION: No difference in Males and Females.