

# Chi-square Test

---

Goodness of Fit and Independence of Attributes

# Chi-square test for independence of attributes

---

*Two random variables  $x$  and  $y$  are called independent if the probability distribution of one variable is not affected by the presence of another. Assume  $O_{ij}$  is the observed frequency count of events belonging to both  $i$ -th category of  $x$  and  $j$ -th category of  $y$ . Also assume  $E_{ij}$  to be the corresponding expected count if  $x$  and  $y$  are independent. The null hypothesis of the independence assumption is to be rejected if the  $p$ -value of the following Chi-squared test statistics is less than a given significance level  $\alpha$ .*

$$\chi^2 = \sum \left[ \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right]$$



*Problem 1 :The below table gives the distribution of students according to the family type and the anxiety level*

| <i>Family type</i>    | <i>Anxiety level</i> |               |             |
|-----------------------|----------------------|---------------|-------------|
|                       | <i>Low</i>           | <i>Normal</i> | <i>High</i> |
| <i>Joint family</i>   | 35                   | 42            | 61          |
| <i>Nuclear family</i> | 48                   | 51            | 68          |

# R- Code and Interpretation

```
> data<-matrix(c(35,42,61,48,51,68),ncol=3,byrow=T)
> data
      [,1] [,2] [,3]
[1,]   35   42   61
[2,]   48   51   68
> chisq.test(data)
```

Pearson's Chi-squared test

```
data: data
X-squared = 0.53441, df = 2, p-value = 0.7655
```

*Here P value (0.7655) > 0.05. Hence there is no evidence to reject the Null hypothesis. So we consider the anxiety level and family type as independent.*



# Problem

---

*In the built-in data set survey, the Smoke column records the students smoking habit, while the Exer column records their exercise level. The allowed values in Smoke are "Heavy", "Regul" (regularly), "Occas" (occasionally) and "Never". As for Exer, they are "Freq" (frequently), "Some" and "None". We can tally the students smoking habit against the exercise level with the table function in R. The result is called the contingency table of the two variables.*

```
> library(MASS)
> tbl = table(survey$Smoke, survey$Exer)
> tbl
```

|       | Freq | None | Some |
|-------|------|------|------|
| Heavy | 7    | 1    | 3    |
| Never | 87   | 18   | 84   |
| Occas | 12   | 3    | 4    |
| Regul | 9    | 1    | 7    |

*Test the hypothesis whether the students smoking habit is independent of their exercise level at .05 significance level.*



# R- Code and Interpretation

---

```
> chisq.test(tbl)
```

```
    Pearson's Chi-squared test
```

```
data:  tbl
```

```
X-squared = 5.4885, df = 6, p-value = 0.4828
```

```
Warning message:
```

```
In chisq.test(tbl) : Chi-squared approximation may be incorrect
```

*As the p-value 0.4828 is greater than the .05 significance level, we do not reject the null hypothesis that the smoking habit is independent of the exercise level of the students.*

# Enhanced Solution

---

*The warning message found in the solution above is due to the small cell values in the contingency table. To avoid such warning, we combine the second and third columns of tbl.*

```
> ctbl = cbind(tbl[, "Freq"], tbl[, "None"] + tbl[, "Some"])
> ct

> ctbl
      [,1] [,2]
Heavy     7    4
Never    87   102
Occas    12    7
Regul     9    8
```



# R- Code

---

```
> chisq.test(ctbl)
```

```
Pearson's Chi-squared test
```

```
data:  ctbl
```

```
X-squared = 3.2328, df = 3, p-value = 0.3571
```

# Goodness of Fit

---

*A biologist is conducting a plant breeding experiment in which plants can have one of four phenotypes. If these phenotypes are caused by a simple Mendelian model, the phenotypes should occur in a 9:3:3:1 ratio. She raises 41 plants with the following phenotypes.*

| <i>Phenotype</i> | <i>1</i> | <i>2</i> | <i>3</i> | <i>4</i> |
|------------------|----------|----------|----------|----------|
| <i>count</i>     | 20       | 10       | 7        | 4        |

*Should she worry that the simple genetic model doesn't work for her phenotypes?*



# R- Code & Inference

```
> plants <- c(20, 10, 7, 4)
> chisq.test(plants, p = c(9/16, 3/16, 3/16, 1/16))

    Chi-squared test for given probabilities

data:  plants
X-squared = 1.9702, df = 3, p-value = 0.5786

Warning message:
In chisq.test(plants, p = c(9/16, 3/16, 3/16, 1/16)) :
  Chi-squared approximation may be incorrect
```

The Chi-squared distribution is only an approximation to the sampling distribution of our test statistic, and the approximation is not very good when the expected cell counts are too small. This is the reason for the warning.

Here the probability value  $p$  is greater than alpha level (0.05), so we do not reject the null hypothesis.

# Fitting of Binomial Distribution with Goodness of Fit

*A survey of 320 families with 5 children each revealed the following distribution:*

|                       |           |           |            |           |           |           |
|-----------------------|-----------|-----------|------------|-----------|-----------|-----------|
| <i>Number of Boys</i> | <i>5</i>  | <i>4</i>  | <i>3</i>   | <i>2</i>  | <i>1</i>  | <i>0</i>  |
| <i>No of Girls</i>    | <i>0</i>  | <i>1</i>  | <i>2</i>   | <i>3</i>  | <i>4</i>  | <i>5</i>  |
| <i>No of families</i> | <i>14</i> | <i>56</i> | <i>110</i> | <i>88</i> | <i>40</i> | <i>12</i> |

*Is this result consistent with the hypothesis that male and female births are equally possible?*

*Solution :*

*Let us setup the null hypothesis that the data are consistent with the hypothesis of equal probability for male and female births.*



# R- CODE & INFERENCE

```
> x=c(5,4,3,2,1,0)
> n=5                                #Probability of 'r' male births in a family
> N=320                              #Total Number of families
> P<-0.5                             #Probability of Male Birth
> Obf<-c(14,56,110,88,40,12)         #Observed frequencies
> exf<-dbinom(x,n,P)*320             #Expected frequencies
> # check the Condition Sum of Observed and Expected are Equal
> sum(Obf)
[1] 320
> sum(exf)
[1] 320
> chisq<-sum((Obf-exf)^2/exf)
> chisq
[1] 7.16
> qchisq(0.95,5)
[1] 11.0705
```

***Calculated value of chi-square is less than the tabulated value ,it is not significant at 5 % level of significance and hence the null hypothesis of equal probability for male and female births.***

# Fitting of Poisson Distribution with Goodness of Fit

---

*Fit a Poisson distribution to the following data and test the goodness of fit*

|     |     |    |    |   |   |   |   |
|-----|-----|----|----|---|---|---|---|
| $X$ | 0   | 1  | 2  | 3 | 4 | 5 | 6 |
| $f$ | 275 | 72 | 30 | 7 | 5 | 2 | 1 |



```
> x<-0:6
> f<-c(275,72,30,7,5,2,1)
> lambda<-(sum(f*x)/sum(f))    #mean
> expf <-dpois(x,lambda)*sum(f)  #expcted frequencies
> f1=round(expf)
> # check obserevd and Expected frequencies Total
> sum(f)
[1] 392
> sum(f1)
[1] 393
> # here subtrat '1' from expected frequencies
> #The last 3 frequencies are less than 5 so ccombine these frequencies in Observation and Expected
> obf<-c(275,72,30,15)
> exf<-c(242,117,28,6)
> chisq<-sum(((obf-exf)^2)/exf)
> chisq
[1] 35.45055
> qchisq(0.95,2)
[1] 5.991465
```

# Inference

---

*Since calculated value of  $\chi^2 = 35.45055$  is much greater than 5.99, it is highly significant. Hence we conclude that poisson distribution is not good fit to the given data*



# Fitting the Normal Distribution with Goodness of Fit

---

*Problem :* The following table displays a frequency distribution of heights of trees in a certain locality. Fit a normal distribution to the data and test the goodness of fit.

| Class Interval | Frequency |
|----------------|-----------|
| 13.20 – 20.90  | 2         |
| 20.90 – 28.60  | 10        |
| 28.60 – 36.30  | 16        |
| 36.30 – 44.00  | 37        |
| 44.00 – 51.70  | 43        |
| 51.70 – 59.40  | 39        |
| 59.40 – 67.10  | 29        |
| 67.10 – 74.80  | 13        |
| 74.80 – 82.50  | 06        |
| 82.50 – 90.20  | 05        |

Heights of Trees (in inches)

```

> midy<-seq(17.05,86.5,length=10)
> f<-c(2,10,16,37,43,39,29,13,6,5)
> mean<-sum(f*midy)/sum(f)
> sd<-sqrt(sum(f*(midy-mean)^2)/sum(f))
> l<-seq(13.2,82.5,length=10)
> l<-c(l,90.2)
> cdf<-pnorm(l,mean,sd)
> cdf<-c(0,cdf,1)
> pcf<-diff(cdf)
> f<-c(0,f,0)
> ex<-round(pcf*sum(f),4)
> fr<-data.frame(f,ex)
> obf<-c(12,16,37,43,39,29,13,11)
> exf<-c(sum(ex[c(1,2,3)]),ex[c(4:9)],sum(ex[c(10,11,12)]))
> sum(obf)
[1] 200
> sum(exf)
[1] 200
> chisq<-sum((obf-exf)^2/exf)
> chisq
[1] 2.153974
> qchisq(0.95,5)
[1] 11.0705

```



# Inference

---

*Here chi-square cal value is less than chi-square tab value then there is no evidence to reject our null hypothesis. ie the fit of normal distribution is good*

# Practice Problems

1. The following data come from a hypothetical survey of 920 people (Men, Women) that ask for their preference of one of the three ice cream flavors (Chocolate, Vanilla, Strawberry). Is there any association between gender and preference for ice cream flavor?

| Gender\flavor | Chocolate | Vanilla | Strawberry |
|---------------|-----------|---------|------------|
| Men           | 100       | 120     | 60         |
| Women         | 350       | 320     | 150        |

2. As a part of quality improvement project focused on a delivery of mail at a department office within a large company, data were gathered on the number of different addresses that had to be changed so that the mail could be redirected to the correct mail stop. Table shows the frequency distribution. Fit binomial distribution and test goodness of fit

| x              | 0 | 1  | 2  | 3  | 4  |
|----------------|---|----|----|----|----|
| f <sub>x</sub> | 5 | 20 | 45 | 20 | 10 |

**The number of Addresses Needing Change**