

## 9. The Chi-Squared Test - Example Problems I

### Chi-squared Test I

1/1得分 (计入成绩)

Let  $\hat{\mathbf{p}}$  denote the MLE for a categorical statistical model  $(\{a_1, \dots, a_K\}, \{\mathbf{P}_{\mathbf{p}}\}_{\mathbf{p} \in \Delta_K})$ . Let  $\mathbf{p}^*$  denote the true parameter. Then  $\sqrt{n}(\hat{\mathbf{p}} - \mathbf{p}^*)$  is asymptotically normal and

$$n \sum_{i=1}^K \frac{(\hat{p}_i - p_i^*)^2}{p_i^*} \xrightarrow[n \rightarrow \infty]{(d)} \chi_{K-1}^2.$$

Consider the particular categorical distribution from the previous problems in this lecture, where we have the statistical experiment  $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{iid}{\sim} \mathbf{P}_{\mathbf{p}}$  and associated statistical model  $(\{1, 2, 3\}, \{\mathbf{P}_{\mathbf{p}}\}_{\mathbf{p} \in \Delta_3})$ . We will use the above fact to hypothesis test between the following null and alternative:

$$\begin{aligned} H_0 : \mathbf{p}^* &= [1/3 \ 1/3 \ 1/3]^T \\ H_1 : \mathbf{p}^* &\neq [1/3 \ 1/3 \ 1/3]^T. \end{aligned}$$

Consider the test

$$\psi = \mathbf{1} \left( n \sum_{i=1}^3 \frac{(\hat{p}_i - \frac{1}{3})^2}{1/3} > C \right),$$

for a threshold  $C$ .

Compute the asymptotic p-value of the test  $\psi$  on the data set

$$\mathbf{x} = 1, 3, 1, 2, 2, 2, 1, 1, 3, 1, 1, 2.$$

Give a numerical value with at least 4 decimals. Use [this tool](#) to find the tail probabilities of a  $\chi^2$  distribution (you may also use any other software). If you are using this tool, note that you need to set "Choose Type of Control" to "Adjust X-axis quantile (Chi square) value" to find the tail probability associated with an x-axis value for a chi-squared distribution with degrees of freedom set in the "Degrees of Freedom" box.

✓ Answer: 0.3679

#### Solution:

Recall from the previous problem on this statistical model that the MLE is given by

$$\hat{\mathbf{p}} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{3} \\ \frac{1}{6} \end{bmatrix}.$$

Therefore, for the 12-point data set  $\mathbf{x}$ ,

$$n \sum_{i=1}^3 \frac{(\widehat{p_i} - \frac{1}{3})^2}{1/3} = 12 \cdot 3 \left( \frac{1}{6^2} + 0 + \frac{1}{6^2} \right) = 2.$$

By the asymptotic normality, we have

$$n \sum_{i=1}^3 \frac{(\widehat{p_i} - \frac{1}{3})^2}{1/3} \xrightarrow[n \rightarrow \infty]{(d)} \chi_2^2.$$

Consulting the link provided, we see that if  $\boldsymbol{X} \sim \chi_2^2$ , then  $P(\boldsymbol{X} \geq 2) \approx \mathbf{36.79\%}$ . Hence, the asymptotic p-value for the test  $\boldsymbol{\psi}$  on this dataset is approximately equal to **0.3679**.

提交

你已经尝试了1次（总共可以尝试3次）

**i** Answers are displayed within the problem

### Playing Dice I

1/1得分 (计入成绩)

You and your friend play dice games for fun, but one day you suspect that the die your friend uses is not fair. You will gather data and use the tools of hypothesis testing in this problem to provide a plausible answer as to whether or not the die is fair.

Your statistical model is  $(\{1, 2, 3, 4, 5, 6\}, \{\mathbf{P}_{\mathbf{p}}\}_{\mathbf{p} \in \Delta_6})$ . You roll the die **15** times, writing the sample as random variables

$\boldsymbol{X}_1, \dots, \boldsymbol{X}_{15} \overset{iid}{\sim} \mathbf{P}_{\mathbf{p}^*}$  where  $\mathbf{p}^*$  is the true parameter. Your null and alternative hypothesis are, respectively,

$$\begin{aligned} H_0 : \mathbf{p}^* &= [1/6 \ 1/6 \ 1/6 \ 1/6 \ 1/6 \ 1/6]^T. \\ H_1 : \mathbf{p}^* &\neq [1/6 \ 1/6 \ 1/6 \ 1/6 \ 1/6 \ 1/6]^T. \end{aligned}$$

Let  $\widehat{\mathbf{p}}$  denote the MLE for the true parameter  $\mathbf{p}^*$ . You use the following test statistic to test the hypotheses:

$$T_n = n \sum_{j=1}^6 \frac{(\hat{p}_j - \frac{1}{6})^2}{\frac{1}{6}}.$$

$T_n$  converges to a random variable  $\boldsymbol{X} \sim \chi_k^2$  for some integer  $k$ . What is  $k$ ?

5

✔ Answer: 5

**Solution:**

We know from this lecture (so far) that if the sample space consists of  $K$  elements, and  $\boldsymbol{X}_1, \dots, \boldsymbol{X}_n \overset{iid}{\sim} \mathbf{P}_{\mathbf{p}^0}$ , then

$$T_n = n \sum_{j=1}^K \frac{(\hat{p}_j - p_j^0)^2}{p_j^0} \xrightarrow[n \rightarrow \infty]{(d)} \chi_{K-1}^2.$$

Our sample space consists of **6** elements (recall that  $\boldsymbol{E} = \{1, 2, 3, 4, 5, 6\}$ ), so we conclude that the limiting distribution is  $\chi_5^2$ .

提交

你已经尝试了1次（总共可以尝试1次）

**i** Answers are displayed within the problem

### Playing Dice II

1/1得分 (计入成绩)

We use the same statistical set-up as above. In particular, our test-statistic is

$$T_n = n \sum_{j=1}^6 \frac{(\hat{p}_j - \frac{1}{6})^2}{\frac{1}{6}}.$$

You use a test of the form

$$\psi_n = \mathbf{1} \left( T_n > C \right).$$

What value of  $C$  should be chosen so that  $\psi$  is a test of asymptotic level **5%**? Give a numerical value with at least 2 decimals.

Use [this table](#) to find the quantiles of a chi-squared distribution.

11.070

✔ Answer: 11.070

Solution:

From the previous problem, we know that

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} \chi_5^2.$$

Consulting a table, we see that if  $X \sim \chi_5^2$ , then if  $C = \mathbf{11.070}$

$$P \left( X > C \right) = 0.05.$$

Therefore,

$$\lim_{n \rightarrow \infty} P_{H_0} \left[ T_n > 11.070 \right] = 0.05.$$

By definition,  $\psi_n$  is a test with asymptotic level **5%**.

Submit

You have used 1 of 3 attempts

❗

Answers are displayed within the problem

Playing Dice III

1/1得分 (计入成绩)

We use the same statistical set-up as above. Recall that you use a test of the form

$$\psi_n = \mathbf{1} \left( T_n > C \right),$$

where  $C$  is a constant chosen in the previous problem so that  $\psi_n$  has asymptotic level **5%**. Suppose you observe that data set

$$\mathbf{5, 6, 1, 6, 4, 1, 2, 4, 6, 6, 1, 6, 6, 3, 5}.$$

Do you **reject** or **fail to reject** the null hypothesis that the die is fair? (You are encouraged to use computational tools.)

☐ Reject

Solution:

The MLE is

$$\hat{\mathbf{p}} = \frac{1}{15} [3 \ 1 \ 1 \ 2 \ 2 \ 6]^T.$$

We compute that for this data set,

$$\begin{aligned} T_{15} &= 15 \left( \frac{(\frac{3}{15} - \frac{1}{6})^2}{1/6} + \frac{(\frac{1}{15} - \frac{1}{6})^2}{1/6} + \frac{(\frac{1}{15} - \frac{1}{6})^2}{1/6} + \frac{(\frac{2}{15} - \frac{1}{6})^2}{1/6} + \frac{(\frac{2}{15} - \frac{1}{6})^2}{1/6} + \frac{(\frac{6}{15} - \frac{1}{6})^2}{1/6} \right) \\ &\approx 7.000. \end{aligned}$$

Since  $C = 11.070$ , by the previous problem,  $\psi$  fails to reject on the given data set.

**Remark:** This is a rather surprising result given that the number 6 has appeared an overwhelming 6 times out of 15 trials and the numbers 2 and 3 have each appeared only once. Without performing this test, one would have probably concluded that the die is likely not a fair die (would have rejected the null hypothesis).

Submit

You have used 1 of 1 attempt

**i** Answers are displayed within the problem

Hypothesis Testing for a Non-Uniform Distribution

1/1得分 (计入成绩)  
Aliens from Planets X, Y, and Z gather on a remote planet every year to decide intergalactic policy. The organizing committee wants to check that the numbers of visitors from each planet is representative of that planet's population. Note that

- Population of Planet X: 1 million
- Population of Planet Y: 4 million
- Population of Planet Z: 5 million.

Let  $E = \{X, Y, Z\}$  denote the sample space. There are a total of 100 visitors chosen for this year's meeting from the overall population of 10 million. Let  $\xi_1, \dots, \xi_{100}$  denote random variables corresponding to alien 1, 2, ..., 100, respectively, so that

$$\xi_i = \begin{cases} X & \text{if alien } i \text{ comes from Planet X} \\ Y & \text{if alien } i \text{ comes from Planet Y} \\ Z & \text{if alien } i \text{ comes from Planet Z} \end{cases}$$

The organizing committee models the outcome of the selection process as a statistical experiment with a categorical distributional model:  $(\{X, Y, Z\}, \{\mathbf{P}_{\mathbf{p}}\}_{\mathbf{p} \in \Delta_3})$  and write  $\xi_1, \dots, \xi_{100} \stackrel{iid}{\sim} \mathbf{P}_{\mathbf{p}^*}$  where  $\mathbf{p}^*$  is the true parameter. The null hypothesis and alternative hypothesis are, respectively,

$$\begin{aligned} H_0 : \mathbf{p}^* &= \begin{bmatrix} \frac{1}{10} \\ \frac{4}{10} \\ \frac{5}{10} \end{bmatrix} . \\ H_1 : \mathbf{p}^* &\neq \begin{bmatrix} \frac{1}{10} \\ \frac{4}{10} \\ \frac{5}{10} \end{bmatrix} . \end{aligned}$$

**Remark:** Note that if  $H_0$  holds, then the visiting delegation is representative of the populations of the three planets in the sense that the percentage of visitors from Planet X (respectively, Planet Y and Z) is not far from the percentage of aliens that live on Planet X (respectively, Planet Y and Z).

Suppose there are 20 visitors from Planet X, 30 visitors from Planet Y, and 50 visitors from Planet Z. Let  $\hat{\mathbf{p}}$  denote the MLE for  $\mathbf{p}^*$  for this data set.

What is the asymptotic p-value of the  $\chi^2$  test

$$\psi_{100} = \mathbf{1} \left( 100 \left( \frac{(\hat{p}_1 - \frac{1}{10})^2}{1/10} + \frac{(\hat{p}_2 - \frac{4}{10})^2}{4/10} + \frac{(\hat{p}_3 - \frac{5}{10})^2}{5/10} \right) > C \right) ?$$

Use [this tool](#) to find the tail probabilities of a  $\chi^2$  distribution (you may also use any software you are familiar with). If you are using this tool, note that you need to set "Choose Type of Control" to "Adjust X-axis quantile (Chi square) value" to find the tail probability associated with an x-axis value for a chi-squared distribution with degrees of freedom set in the "Degrees of Freedom" box.

Give a numerical value with at least 5 decimals. (You are encouraged to also use computational tools.)

0.00193

✔ Answer: 0.00193

Solution:

Since the sample space has 3 elements, X, Y, and Z, we know that

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} \chi^2_2.$$

Recall that the asymptotic p-value is the smallest asymptotic level so that the test  $\psi_n$  rejects on the given data set.

The maximum likelihood estimator is given by

$$\hat{\mathbf{p}} = \begin{bmatrix} \frac{2}{10} \\ \frac{3}{10} \\ \frac{5}{10} \end{bmatrix}.$$

Thus the test statistic takes the value

$$T_{100} = 100 \left( \frac{(\frac{2}{10} - \frac{1}{10})^2}{1/10} + \frac{(\frac{3}{10} - \frac{4}{10})^2}{4/10} + \frac{(\frac{5}{10} - \frac{5}{10})^2}{5/10} \right) \approx 12.5.$$

Let  $X \sim \chi^2_2$ . To compute the asymptotic p-value, we set the threshold  $C$  to be equal to the observed value of  $T_{15}$  and compute

$$P(X \geq T_{100}) = P(X \geq 12.5) \approx 0.00193$$

using the online tool.

Submit

You have used 1 of 3 attempts

📘 Answers are displayed within the problem