

<u>Lecture 15: Goodness of Fit Test for</u> 11. Chi-Squared Test for a Family of

> Discrete Distributions

<u>Course</u> > <u>Unit 4 Hypothesis testing</u> > <u>Discrete Distributions</u>

11. Chi-Squared Test for a Family of Discrete Distributions

In the problems on this page, you will apply the χ^2 goodness of fit test to determine whether or not a sample has a binomial distribution.

So far, we have used the χ^2 test to determine if our data had a categorical distribution with specific parameters (e.g. uniform on an Nelement set).

For the problems on this page, we extend the discussion on χ^2 tests **beyond** what was discussed in lecture to the following more general statistical set-up.

Let $X_1,\ldots,X_n\stackrel{iid}{\sim} X\sim {f P}$ denote iid discrete random variables supported on $\{0,\ldots,K\}$. We will decide between the following null and alternative hypotheses:

 $H_0: \quad \mathbf{P} \in \left\{ \mathrm{Bin}\left(K, heta
ight)
ight\}_{ heta \in (0,1)}$

 $H_1: \mathbf{P} \notin \{\operatorname{Bin}(K,\theta)\}_{\theta \in (0,1)},$

where the null hypothesis can be rephrased as:

 $H_0: \quad ext{there exists } heta \in (0,1) ext{ such that for all } j=0,\ldots,K, ext{ we have } P\left(X=j
ight) = {K \choose j} heta^j (1- heta)^{K-j}.$

Review: Log-likelihood for a Binomial Distribution

2/2 points (graded)

Let $(\{0,\ldots,K\},\{\operatorname{Bin}\,(K, heta)\}_{ heta\in(0,1)})$ denote a binomial statistical model. Let $X_1,\ldots,X_n\stackrel{iid}{\sim}\operatorname{Bin}\,(K, heta^*)$ for some unknown parameter $\theta^* \in (0,1)$.

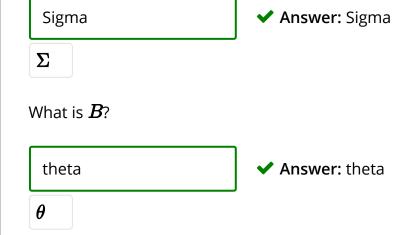
The log-likelihood of this statistical model can be written

$$C + A \log B + (nK - A) \log (1 - B)$$

where C is independent of heta, A depends on $\sum_{i=1}^n X_i$, and B depends on heta.

What is A?

Use **Sigma** to stand for $\sum_{i=1}^n X_i$.



STANDARD NOTATION

Solution:

The pmf of $\operatorname{Bin}\left(K,\theta\right)$ is

$$j\mapsto inom{K}{j} heta^j(1- heta)^{K-j}$$

for $j \in \{1, \dots, K\}$.

Therefore, the likelihood is given by

$$egin{aligned} L_n\left(X_1,\ldots,X_n, heta
ight) &= \prod_{i=1}^n \left(inom{K}{X_i} heta^{X_i} (1- heta)^{K-X_i}
ight) \ &= \left(\prod_{i=1}^n inom{K}{X_i}
ight) heta^{\sum_{i=1}^n X_i} (1- heta)^{nK-\sum_{i=1}^n X_i}. \end{aligned}$$

Taking the logarithm, we have

$$\log L_n\left(X_1,\ldots,X_n, heta
ight) = \log \left(\prod_{i=1}^n {K\choose X_i}
ight) + \left(\sum_{i=1}^n X_i
ight) \log heta + \left(nK - \sum_{i=1}^n X_i
ight) \log \left(1- heta
ight).$$

Therefore, $A = \sum_{i=1}^n X_i$ and B = heta.

Submit

You have used 1 of 4 attempts

1 Answers are displayed within the problem

Review: MLE for a Binomial Distribution

1/1 point (graded)

As above, let $(\{0,\ldots,K\},\{\operatorname{Bin}\,(K,\theta)\}_{\theta\in(0,1)})$ denote a binomial statistical model. Let $X_1,\ldots,X_n\stackrel{iid}{\sim}\operatorname{Bin}\,(K,\theta^*)$ for some unknown parameter $\theta^*\in(0,1)$.

Which of the following denotes the MLE for θ^* ?

- $\circ \sum_{i=1}^n X_i$
- $\bigcup \frac{1}{n} \sum_{i=1}^{n} X_i$
- $0 \frac{1}{K} \sum_{i=1}^{n} X_i$
- \bullet $\frac{1}{nK} \sum_{i=1}^{n} X_i \checkmark$

Solution:

Recall from the previous problem that

$$\log L_n\left(X_1,\ldots,X_n, heta
ight) = C + \left(\sum_{i=1}^n X_i
ight) \log heta + \left(nK - \sum_{i=1}^n X_i
ight) \log \left(1 - heta
ight)$$

where $oldsymbol{C}$ does not depend on $oldsymbol{ heta}$.

To compute the MLE, we need to maximize the above with respect to the parameter heta. We set the derivative to be 0:

$$0 = \frac{\sum_{i=1}^n X_i}{\theta} - \frac{nK - \sum_{i=1}^n X_i}{1 - \theta}.$$

The above holds when

$$p = rac{1}{nK} \sum_{i=1}^n X_i.$$

Therefore, the right-hand side is the MLE for this statistical model.

Submit

You have used 1 of 2 attempts

• Answers are displayed within the problem

χ^2 -Test for a Family of Distributions :

Now, we return to the following more general statistical set-up.

Let $X_1, \ldots, X_n \overset{iid}{\sim} \mathbf{P}$ denote iid discrete random variables supported on $\{0, \ldots, K\}$. We will decide between the following null and alternative hypotheses.

 $H_0: \ \ \mathbf{P} \in \left\{ \mathrm{Bin}\left(K, heta
ight)
ight\}_{ heta \in (0,1)}.$

 $H_1: \quad \mathbf{P}
otin \left\{ \mathrm{Bin} \left(K, heta
ight)
ight\}_{ heta \in (0,1)}.$

Let $f_{ heta}$ denote the pmf of the distribution $\mathrm{Bin}\,(K, heta)$, and let $\hat{ heta}$ denote the MLE of the parameter heta from the previous problem.

Further, let N_j denote the number of times that j ($j\in\{0,1,\ldots,K\}$) appears in the data set X_1,\ldots,X_n (so that $\sum_{j=0}^K N_j=n$.) The

 χ^2 test statistic for this hypothesis test is defined to be

$$T_n := n \sum_{j=0}^K rac{\left(rac{N_j}{n} - f_{\hat{ heta}}\left(j
ight)
ight)^2}{f_{\hat{ heta}}\left(j
ight)}.$$

This statistic is different from before. Previously, under the null hypothesis, $\mathbf{P}(X=j)=p_j$ for some fixed p_j . Here, instead, we use $f_{\hat{\theta}}(j)$ to estimate $\mathbf{P}(X=j)$. This statistic still converges in distribution to a χ^2 distribution, but the number of degrees of freedom is smaller.

Degrees of Freedom for χ^2 Test for a Family of Distribution

More generally, to test if a distribution \mathbf{P} is described by some member of a family of discrete distributions $\{\mathbf{P}_{\theta}\}_{\theta\in\Theta\subset\mathbb{R}^d}$ where $\Theta\subset\mathbb{R}^d$ is d-dimensional, with support $\{0,1,2,\ldots,K\}$ and pmf f_{θ} , i.e. to test the hypotheses:

 $H_0: \ \ \mathbf{P} \in \{\mathbf{P}_{ heta}\}_{ heta \in \Theta}$

 $H_1: \mathbf{P} \notin \{\mathbf{P}_{\theta}\}_{\theta \in \Theta},$

then if indeed $\mathbf{P} \in \{\mathbf{P}_{\theta}\}_{\theta \in \Theta \subset \mathbb{R}^d}$ (i.e., the null hypothesis H_0 holds), and if in addition some technical assumptions hold, then we have that

$$T_n := n \sum_{j=0}^K rac{\left(rac{N_j}{n} - f_{\hat{ heta}}\left(j
ight)
ight)^2}{f_{\hat{ heta}}\left(j
ight)} \stackrel{(d)}{\longrightarrow} \chi^2_{(K+1)-d-1}.$$

Note that K+1 is the support size of $\mathbf{P}_{ heta}$ (for all heta.)

In our example testing for a binomial distribution, the parameter θ is one-dimensional, i.e. d=1. Therefore, under the null hypothesis H_0 , it holds that

$$T_n \xrightarrow[n o \infty]{(d)} \chi^2_{(K+1)-1-1} = \chi^2_{K-1}.$$

Chi-squared Test for a Binomial Distribution on a Sample Data Set I

1/1 point (graded)

Consider the same statistical set-up as above. In particular, we have the test statistic

$$T_n := n \sum_{j=0}^K rac{\left(rac{N_j}{n} - f_{\hat{ heta}}\left(j
ight)
ight)^2}{f_{\hat{ heta}}\left(j
ight)}.$$

where $\hat{ heta}$ is the MLE for the binomial statistical model $(\{0,1,\ldots,K\},\{\operatorname{Bin}\,(K, heta)\}_{ heta\in(0,1)})$.

We define our test to be

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

where au is a threshold that you will specify. For the remainder of this page, we will assume that K=3 (the sample space is $\{0,1,2,3\}$).

What value of au should be chosen so that ψ_n is a test of asymptotic level 5%? Give a numerical value with at least 3 decimals.

(Use this table or software to find the quantiles of a chi-squared distribution.)

$$\tau = \begin{bmatrix} 5.991 & \checkmark \text{ Answer: } 5.991 \end{bmatrix}$$

Solution:

Since K=3 and d=1, we know that the limiting distribution of T_n is χ^2_2 . Thus, the asymptotic level is the value au such that

$$\lim_{n o\infty}P\left(T_{n}> au
ight)=P\left(Z> au
ight)=0.05$$

where $Z\sim\chi^2_2$. Hence, au should be chosen to be 5.991 (from the given table).

Submit

You have used 1 of 2 attempts

Answers are displayed within the problem

Chi-squared Test for a Binomial Distribution on a Sample Data Set II

3/3 points (graded)

Consider the same statistical set-up as above. Suppose we observe a data set consisting of 1000 observations as described in the following (format: i, number of observations of i):

- $i N_i$
- 0 339
- 1 455
- 2 180
- 3 26

What is the value of the test statistic T_n for this data set? Give a numerical value with at least 4 decimals. (You are encouraged to use computational software.)

$$T_n = \begin{bmatrix} 0.88286 \end{bmatrix}$$
 Answer: 0.8829

What is the p-value of this data set with respect to the test ψ_{1000} ? Give a numerical value with at least 4 decimals.

Use <u>this tool</u> to find the tail probabilities of a χ^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.

If ψ_n is designed to have level 5%, would you **reject** or **fail to reject** on the given data set?

Reject

Fail to reject

Solution:

Observe that the MLE is given by

$$\hat{p} = rac{1}{3 \cdot 1000} (455 + 2 \cdot 180 + 3 \cdot 26) pprox 0.29767.$$

Thus for this data set,

$$T_n = 1000 \cdot \left(rac{\left(rac{339}{1000} - \left(rac{3}{0}
ight) \left(0.2977^0
ight) \left(0.7023
ight)^{3-0}
ight)^2}{\left(rac{3}{0}
ight) \left(0.2977^0
ight) \left(0.7023
ight)^{3-0}} + rac{\left(rac{455}{1000} - \left(rac{3}{1}
ight) \left(0.2977^1
ight) \left(0.7023
ight)^{3-1}
ight)^2}{\left(rac{3}{1000} - \left(rac{3}{2}
ight) \left(0.2977^2
ight) \left(0.7023
ight)^{3-2}
ight)^2} + rac{\left(rac{26}{1000} - \left(rac{3}{3}
ight) \left(0.2977^3
ight) \left(0.7023
ight)^{3-3}
ight)^2}{\left(rac{3}{2}
ight) \left(0.2977^2
ight) \left(0.7023
ight)^{3-2}} + rac{\left(rac{26}{1000} - \left(rac{3}{3}
ight) \left(0.2977^3
ight) \left(0.7023
ight)^{3-3}
ight)^2}{\left(rac{3}{3}
ight) \left(0.2977^3
ight) \left(0.7023
ight)^{3-3}}
ight)} pprox 0.8829$$

The asymptotic p-value for this data set is given by

$$\lim_{n o\infty}P\left(T_{n}>0.8829
ight)=P\left(Z>0.8829
ight).$$

where $Z\sim\chi^2_2$. Consulting the suggested link, we see that $P\left(Z>0.8829
ight)pprox0.6431$.

According to the golden rule of p-values, since 0.6431 > 0.05, we should **fail to reject** the null hypothesis that X_1, \ldots, X_{1000} are distributed as Bin(3, p) for some value of the parameter p.

Submit

You have used 2 of 3 attempts

Answers are displayed within the problem

Discussion

Show Discussion

Topic: Unit 4 Hypothesis testing:Lecture 15: Goodness of Fit Test for Discrete Distributions / 11. Chi-Squared Test for a Family of Discrete Distributions