



Lecture 15: Goodness of Fit Test for

5. Maximum Likelihood Estimator

Course > Unit 4 Hypothesis testing > Discrete Distributions

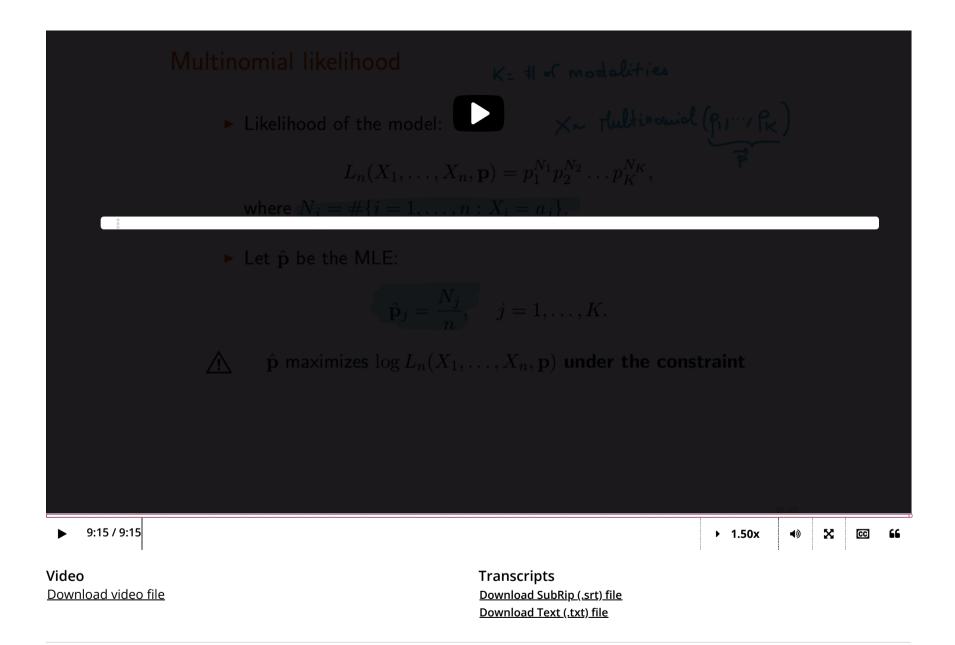
> for the Categorical Distribution

5. Maximum Likelihood Estimator for the Categorical Distribution

Video note: The following video derives the maximum likelihood estimator $\hat{\mathbf{p}}$ of a categorical statistical model. Note that we have seen this previously in Lecture 10 and in Recitation 6.

The maximum likelihood estimator will form the basis of goodness of fit testing for discrete distributions.

MLE for the Categorical Distribution



Concept Check: Examples of the Categorical Distribution

2/2 points (graded)

Consider the distribution $\mathrm{Ber}\,(0.25)$. Consider the categorical statistical model $(\{a_1,\ldots,a_K\},\{\mathbf{P_p}\})$ for this Bernoulli distribution.

If we let $a_1=1$ and $a_2=0$, then this corresponds to a categorical distribution ${f P_p}$ with parameter vector ${f p}$ given by...

0.25

 $\bigcirc 0.75$

 $\bullet \begin{bmatrix} 0.25 \;\; 0.75 \end{bmatrix}^T$

 $\bigcirc [0.75 \ 0.25]^T$



Let $a_i=i$ for $i=1,\ldots,K$. The uniform distribution on $E=\{1,2,\ldots,K\}$ can be expressed as a categorical distribution $\mathbf{P}_{\mathbf{p}}$ for some choice of parameter \mathbf{p} .

What is $\sum_{i=1}^K p_i^2$?

1/K

✓ Answer: 1/K

 $\frac{1}{K}$

STANDARD NOTATION

Solution:

Let $X \sim \mathrm{Ber}\,(0.25)$. Observe that

$$p_1 = P(X = a_1) = P(X = 1) = 0.25$$

and

$$p_2 = P(X = a_2) = P(X = 0) = 0.75.$$

Hence, $\mathbf{p} = [0.25 \ 0.75]^T$.

Remark: Observe that Ber(p) has a one-dimensional parameter and P_p for this example involves a parameter that is two-dimensional, but such that the second parameter depends on the first one ($p_1=1-p_2$). In general, the categorical distribution for $\mathbf{p}\in\Delta_K$ involves a Kdimensional parameter, but the last parameter p_K , for example, is redundant because $p_K = 1 - \sum_{i=1}^{K-1} p_i$. Even though \mathbf{p} is K-dimensional, the categorical distribution has only K-1 degrees of freedom. This will make our analysis more challenging: the extra constraint on the parameter $\sum_{i=1}^{K} p_i = 1$ implies that the Fisher information for the model as specified **does not exist**. Hence, we cannot apply Wald's test

For the second question, by definition, the uniform distribution weighs all elements in $\{1,\ldots,K\}$ equally. Let **P** denote the parameter vector of the uniform distribution on $\{1, 2, \dots, K\}$. Then

$$p_i = P\left(X=i
ight) = rac{1}{K}.$$

Thus,

$$\sum_{i=1}^K p_i^2 = \sum_{i=1}^K rac{1}{K^2} = rac{1}{K}.$$

Submit

You have used 1 of 2 attempts

1 Answers are displayed within the problem

Likelihood for a Categorical Distribution

3/3 points (graded)

Suppose that K=3, and let $E=\{1,2,3\}$. Let $X_1,\ldots,X_n\stackrel{iid}{\sim}\mathbf{P_p}$ for some unknown $\mathbf{p}\in\Delta_3$. Let $f_\mathbf{p}$ denote the pmf of $\mathbf{P_p}$ and recall that the likelihood is defined to be

$$L_{n}\left(X_{1},\ldots,X_{n},\mathbf{p}
ight)=\prod_{i=1}^{n}f_{\mathbf{p}}\left(X_{i}
ight).$$

Here we let the sample size be n=12, and you observe the sample $\mathbf{x}=x_1,\dots,x_{12}$ given by

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

The likelihood for this data set can be expressed as $L_{12}(\mathbf{x}, \mathbf{p}) = p_1^A p_2^B p_3^C$.

Fill in the values of A, B, and C below.

Solution:

Since K=3 and $E=\{1,2,3\}$,

$$f_{\mathbf{p}}\left(i
ight)=p_{i},\quad i=1,2,3.$$

Next,

$$L_{n}\left(X_{1},\ldots,X_{n},\mathbf{p}
ight)=\prod_{i=1}^{n}f_{\mathbf{p}}\left(X_{i}
ight)=p_{1}^{N_{1}}p_{2}^{N_{2}}p_{3}^{N_{3}}$$

where

$$N_i = ext{number of times } i ext{ appears in } (X_1, \ldots, X_n) \,, \quad i = 1, 2, 3.$$

In the data set above, 1 appears 6 times, 2 appears 4 times, and 3 appears 2 times. Thus, $A=N_1=6$, $B=N_2=4$, and $C=N_3=2$.

Submit

You have used 1 of 2 attempts

• Answers are displayed within the problem

Maximum Likelihood Estimator for Categorical Distribution

3/3 points (graded)

As above, under the statistical model $(\{1,2,3\},\{\mathbf{P_p}\}_{\mathbf{p}\in\Delta_3})$, we have

$$L_{12}\left(\mathbf{x},\mathbf{p}
ight)=p_{1}^{A}p_{2}^{B}p_{3}^{C}$$

where

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

In the previous problem, you found the specific values for A, B, and C.

Recall that the MLE is given by

$$\widehat{\mathbf{p}}_{n}^{MLE} = \operatorname{argmax}_{\mathbf{p} \in \Delta_{3}} \log L_{n}\left(X_{1}, \ldots, X_{n}, \mathbf{p}
ight).$$

By the theory of Lagrange multipliers, one can show that the maximum occurs at the point ${f p}$ such that there exists $\lambda
eq 0$ so that

$$abla \log L_n\left(X_1,\ldots,X_n,\mathbf{p}
ight) = \lambda egin{bmatrix} 1 \ 1 \ 1 \end{bmatrix}.$$

(The gradient above is taken with respect to the parameter **p**.)

Using this result and the previous problem, what is the estimate $\widehat{\mathbf{p}}_{12}^{MLE}$ for \mathbf{p} given the data set \mathbf{x} ?

$$(\widehat{\mathbf{p}}_{12}^{MLE})_1 = \boxed{1/2}$$
 \checkmark Answer: 1/2 $(\widehat{\mathbf{p}}_{12}^{MLE})_2 = \boxed{1/3}$ \checkmark Answer: 1/3 $(\widehat{\mathbf{p}}_{12}^{MLE})_3 = \boxed{1/6}$

Answer: 1/6

Solution:

In the previous problem, we saw that A=6, B=4, and C=2. Thus

$$\log L_n\left(\mathbf{x},\mathbf{p}
ight) = 6\log p_1 + 4\log p_2 + 2\log p_3.$$

Hence,

$$abla \log L_n\left(\mathbf{x},\mathbf{p}
ight) = egin{bmatrix} rac{6}{p_1} \ rac{4}{p_2} \ rac{2}{p_3} \end{bmatrix}.$$

Applying the Lagrange multipliers, we have

$$egin{bmatrix} rac{6}{p_1} \ rac{4}{p_2} \ rac{2}{p_3} \end{bmatrix} = \lambda egin{bmatrix} 1 \ 1 \ 1 \end{bmatrix}.$$

Therefore,

$$p_1=rac{6}{\lambda},\; p_2=rac{4}{\lambda},\; p_3=rac{2}{\lambda}.$$

By the constraint $p_1+p_2+p_3=1$, we see that

$$\lambda = 6 + 4 + 2 = 12.$$

Therefore,

$$\widehat{\mathbf{p}}_{12}^{MLE} = egin{bmatrix} rac{1}{2} \ rac{1}{3} \ rac{1}{6} \end{bmatrix}.$$

Submit

You have used 1 of 2 attempts

1 Answers are displayed within the problem

Discussion

Hide Discussion

	Add a Post
Show all posts ▼	by recent activity ▼
There are no posts in this topic yet.	
×	

© All Rights Reserved