

Topic 21
Goodness of Fit
Fit of a Distribution

Outline

Fit of a Distribution

Blood Bank

Likelihood Function

Likelihood Ratio

Lagrange Multipliers

Hanging Chi-Gram

Fit of a Distribution

Goodness of fit tests examine the case of a sequence of independent observations each of which can have 1 of k possible categories. For example, each of us has one of 4 possible of **blood types**, O , A , B , and AB . The local blood bank has good information from a national database of the **fraction of individuals** having each blood type,

$$\pi_O, \pi_A, \pi_B, \text{ and } \pi_{AB}.$$

The **actual fraction** p_O, p_A, p_B , and p_{AB} of these blood types in the community for a given blood bank may be different than what is seen in the national database. As a consequence, the local blood bank may choose to alter its distribution of blood supply to more accurately reflect local conditions.

Introduction

To place this assessment strategy in terms of formal hypothesis testing, let $\pi = (\pi_1, \dots, \pi_k)$ be postulated values of the probability

$$P_\pi\{\text{individual is a member of } i\text{-th category}\} = \pi_i$$

and let $\mathbf{p} = (p_1, \dots, p_k)$ denote the possible **states of nature**. Then, the **parameter space** is

$$\Theta = \{\mathbf{p} = (p_1, \dots, p_k); p_i \geq 0 \text{ for all } i = 1, \dots, k, \sum_{i=1}^k p_i = 1\}.$$

This parameter space has $k - 1$ **free parameters**. Once these are chosen, the remaining parameter value is determined by the requirement that the sum of the p_i equals 1.

Thus, $\dim(\Theta) = k - 1$.

Overview

- The **hypothesis** is

$$H_0 : p_i = \pi_i, \text{ for all } i = 1, \dots, k \quad \text{versus} \quad H_1 : p_i \neq \pi_i, \text{ for some } i = 1, \dots, k$$

- The **parameter space** for the null hypothesis is a **single point** $\pi = (\pi_1, \dots, \pi_k)$. Thus, $\dim(\Theta_0) = 0$.
- Consequently, the **likelihood ratio test statistic** has **chi-square distribution** with $\dim(\Theta) - \dim(\Theta_0) = k - 1$ **degrees of freedom**.
- The **data** $\mathbf{x} = (x_1, \dots, x_n)$ are the **categories** for each of the n observations.

Likelihood Function

Let's use the **likelihood ratio criterion** to create a test for the distribution of human blood types in a given population. For the **data**

$$\mathbf{x} = \{O, B, O, A, A, A, A, A, O, AB\}$$

in the case of independent observations, the likelihood is

$$L(\mathbf{p}|\mathbf{x}) = p_O \cdot p_B \cdot p_O \cdot p_A \cdot p_A \cdot p_A \cdot p_A \cdot p_A \cdot p_O \cdot p_{AB} = p_O^3 p_A^5 p_B p_{AB}.$$

Notice that the likelihood has a factor of p_i whenever an observation take on the value i . In other words, if we summarize the data using

$$n_i = \#\{\text{observations from category } i\}$$

to create $\mathbf{n} = (n_1, n_2, \dots, n_k)$, a vector that records the number of observations in each category, then, the **likelihood function**

$$L(\mathbf{p}|\mathbf{n}) = p_1^{n_1} \cdots p_k^{n_k}.$$

Likelihood Ratio

The **likelihood ratio** is the ratio of the maximum value of the likelihood under the null hypothesis and the maximum likelihood for any parameter value. In this case, the numerator is the likelihood evaluated at π .

$$\Lambda(\mathbf{n}) = \frac{L(\pi|\mathbf{n})}{L(\hat{\mathbf{p}}|\mathbf{n})} = \frac{\pi_1^{n_1} \pi_2^{n_2} \cdots \pi_k^{n_k}}{\hat{p}_1^{n_1} \hat{p}_2^{n_2} \cdots \hat{p}_k^{n_k}} = \left(\frac{\pi_1}{\hat{p}_1}\right)^{n_1} \cdots \left(\frac{\pi_k}{\hat{p}_k}\right)^{n_k}$$

To find the maximum likelihood estimator $\hat{\mathbf{p}}$, we, as usual, begin by taking the logarithm of the likelihood,

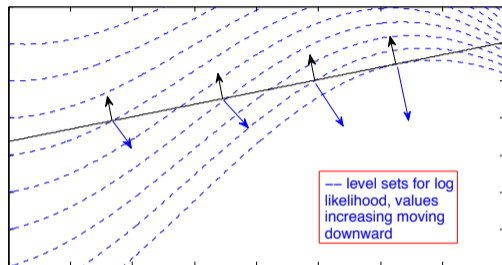
$$\ln L(\mathbf{p}|\mathbf{n}) = \sum_{i=1}^k n_i \ln p_i.$$

Not every set of values for p_i is admissible. So, we cannot just take derivatives, set them equal to 0 and solve. Indeed, we must find a maximum under the **constraint**

$$s(\mathbf{p}) = \sum_{i=1}^k p_i = 1.$$

Lagrange Multipliers

- **Level sets** of the log-likelihood function shown in dashed blue.
- The **constraint** is level set $\{s(\mathbf{p}) = 1\}$ shown in black.
- The **gradients** are indicated by arrows.
- At the maximum, these two arrows are parallel. Their ratio λ is called the **Lagrange multiplier**.



$$\begin{aligned} \nabla_{\mathbf{p}} \ln L(\hat{\mathbf{p}}|\mathbf{n}) &= \lambda \nabla_{\mathbf{p}} s(\hat{\mathbf{p}}). \\ \left(\frac{\partial}{\partial p_1} \ln L(\hat{\mathbf{p}}|\mathbf{n}), \dots, \frac{\partial}{\partial p_k} \ln L(\hat{\mathbf{p}}|\mathbf{n}) \right) &= \lambda \left(\frac{\partial}{\partial p_1} s(\mathbf{p}), \dots, \frac{\partial}{\partial p_k} s(\mathbf{p}) \right) \\ \left(\frac{n_1}{\hat{p}_1}, \dots, \frac{n_k}{\hat{p}_k} \right) &= \lambda (1, \dots, 1) \end{aligned}$$

Lagrange Multipliers

$$\frac{n_i}{\hat{p}_i} = \lambda, \quad n_i = \lambda \hat{p}_i \quad \text{for all } i = 1, \dots, k.$$

Now sum this equality for all values of i and use the **constraint** $s(\mathbf{p}) = 1$ to obtain

$$n = \sum_{i=1}^k n_i = \lambda \sum_{i=1}^k \hat{p}_i = \lambda s(\hat{\mathbf{p}}) = \lambda.$$

Thus, we have that

$$\frac{n_i}{\hat{p}_i} = n \quad \text{and} \quad \hat{p}_i = \frac{n_i}{n}.$$

The estimate for p_i is the fraction of observations in category i . Thus, for the blood bank example,

$$\hat{p}_O = \frac{3}{10}, \quad \hat{p}_A = \frac{5}{10}, \quad \hat{p}_B = \frac{1}{10}, \quad \text{and} \quad \hat{p}_{AB} = \frac{1}{10}.$$

Likelihood Ratio

Next, we substitute the **maximum likelihood estimates** $\hat{p}_i = n_i/n$ into the likelihood ratio to obtain

$$\Lambda(\mathbf{n}) = \frac{L(\boldsymbol{\pi}|\mathbf{n})}{L(\hat{\mathbf{p}}|\mathbf{n})} = \left(\frac{\pi_1}{n_1/n}\right)^{n_1} \cdots \left(\frac{\pi_k}{n_k/n}\right)^{n_k} = \left(\frac{n\pi_1}{n_1}\right)^{n_1} \cdots \left(\frac{n\pi_k}{n_k}\right)^{n_k}.$$

Let $\mathbf{N} = (N_1, \dots, N_k)$ denote the random vector of **observed number of occurrences** for each category i . When the **null hypothesis** holds true,

$$-2 \ln \Lambda(\mathbf{N}) = -2 \sum_{i=1}^k N_i \ln \frac{n\pi_i}{N_i} = 2 \sum_{i=1}^k N_i \ln \frac{N_i}{n\pi_i}$$

has approximately a χ_{k-1}^2 distribution.

Likelihood Ratio

Using the notation $O_i = n_i$ for the number of **observed** occurrences of i and $E_i = n\pi_i$ for the number of **expected** occurrences of i as given by H_0 , we can write the test statistic as

$$G^2 = -2 \ln \Lambda_n(\mathbf{O}) = 2 \sum_{i=1}^k O_i \ln \frac{O_i}{E_i}.$$

The traditional method for a test of goodness of fit, we use, instead of the G^2 **statistic**, an approximation

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

In either case the **p-value** will be the **probability** that the a χ_{k-1}^2 random variable takes a value greater than the test statistic.

Chi-Square Statistic

The **Red Cross** recommends that a blood bank maintains 44% type **O**, 42% type **A**, 10% type **B**, and 4% type **AB**. You suspect that the distribution of blood types in Tucson is not the same as the recommendation. In this case, the hypothesis is

$$H_0 : p_O = 0.44, p_A = 0.42, p_B = 0.10, p_{AB} = 0.04$$

versus

H_1 : at least one p_i is unequal to the given values.

Based on 400 observations, we observe 228 for type **O**, 124 for type **A**, 40 for type **B** and 8 for type **AB**. We find the expected occurrences by computing $400 \times p_i$ using the values in H_0 . This gives the table

type	O	A	B	AB
observed	228	124	40	8
expected	176	168	40	16

Blood Bank

Enter the observations and the proportions under H_0 in the `chisq.test` command. R computes the expected number of observations.

```
> chisq.test(c(228,124,40,8),p=c(0.44,0.42,0.10,0.04))
```

```
Chi-squared test for given probabilities
```

```
data:  c(228, 124, 40, 8)
```

```
X-squared = 30.8874, df = 3, p-value = 8.977e-07
```

The number of degrees of freedom is $4 - 1 = 3$. Note that the p -value is very low and so the distribution of blood types in Tucson is **very unlikely** to be the same as the national distribution.

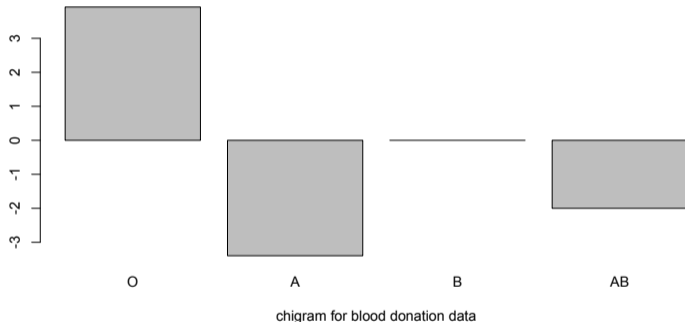
Exercise. Compute by hand the χ^2 statistic from the blood bank data. Use the `pchisq` command to determine the p -value.

Hanging Chi-Gram

To visualize the discrepancies from the null hypothesis, we use a **hanging chi-gram**. This plots category i with a bar, height

$$\frac{O_i - E_i}{\sqrt{E_i}}$$

Note that these values can be either **positive** or **negative**.



```
> resid<-(O-E)/sqrt(E)
> barplot(resid, names.arg=c("O","A","B","AB"),
  xlab="chigram for blood donation data")
```