

Chi-Square ~~Tests of Fit~~ Tests

Suppose I walk around asking people to tell me a random number from 1 to 10.
(integer)

I am interested in whether people are more likely to tell me one number than others. I can phrase this as a hypothesis testing problem:

$$H_0 : P_1 = P_2 = \dots = P_{10} = \frac{1}{10}$$

$$H_1 : \text{otherwise}$$

where $P_i = \text{Prob}(i)$.

More generally, suppose an experiment has k outcomes. Let X denote the outcome. $(1, 2, \dots, 10)$. Set

$$P_i = P(X=i).$$

We are interested in testing

$$H_0: P_i = P_{i0} \text{ for } i=1, \dots, k$$

vs.

$$H_1: \text{otherwise}$$

where $P_{10}, P_{20}, \dots, P_{k0}$ are fixed, known, and satisfy

$$\textcircled{1} \quad 0 < P_{i0} < 1 \quad \text{for all } i$$

$$\textcircled{2} \quad \sum_{i=1}^k P_{i0} = 1.$$

In the example, $k = 10$ and $p_{10} = \frac{1}{10}$
for each i .

What can we do? Well, if $k = 2$ we
know how to proceed:

Let X_1, \dots, X_n be a random sample,
and let Y_1 denote the number of
times $X_j = 1$ occurs in the sample.

Then $Y_1 \sim \text{binom}(n, p_1)$.

By the CLT, $\frac{Y_1 - np_{10}}{\sqrt{np_{10}(1-p_{10})}} \approx N(0,1)$.

if H_0 is true.

Therefore, if we accept ~~H₀~~ H_0 iff

$$\frac{|Y_1 - np_{10}|}{\sqrt{np_{10}(1-p_{10})}} \leq z_{\alpha/2},$$

we have a test w/ level of significance α .
approximately

To generalize this to $k > 2$, it is helpful to consider a different test, also based on the CLT:

$$\text{Let } Z = \frac{Y_1 - np_1}{\sqrt{np_1(1-p_1)}}$$

$$\text{and } Q_1 = Z^2.$$

$$\text{Then } Q_1 \approx \chi^2(1).$$

Thus another test of H_0 vs. H_1 is

to accept H_0 iff

$$Q_1 \leq \chi_{\alpha}^2(1).$$

We may rewrite Q_1 as follows:

$$\begin{aligned} Q_1 &= Z^2 = \frac{(Y_1 - np_1)^2}{np_1(1-p_1)} \\ &= \frac{(Y_1 - np_1)^2}{np_1} + \frac{(Y_1 - np_1)^2}{n(1-p_1)}. \end{aligned}$$

If is true that

$$\bullet \quad p_2 = 1 - p_1$$

$$\bullet \quad (Y_2 - np_2)^2 = (Y_1 - np_1)^2$$

$$\boxed{Y_2 = n - Y_1}$$

simple
algebra

Therefore

$$Q_1 = \frac{(Y_1 - np_1)^2}{np_1} + \frac{(Y_2 - np_2)^2}{np_2}$$

Return to the general case :

- X has k outcomes
- X_1, \dots, X_n is a random sample
- $Y_k = \#\{X_j = k\}$
- $P_i = P(X = i)$

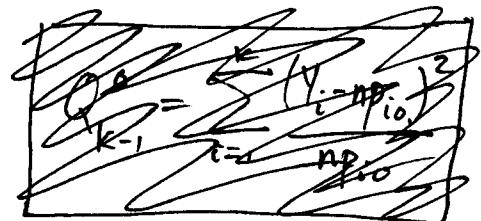
FACT If $Q_{k-1} = \sum_{i=1}^k \frac{(Y_i - np_i)^2}{np_i}$

then $Q_{k-1} \approx \chi^2(k-1)$.

(we will not prove this).

Therefore, to test H_0 vs. H_1 at a level of significance α , we accept H_0 when

$$Q_{k-1} = \sum_{i=1}^k \frac{(Y_i - np_{io})^2}{np_{io}} \leq \chi_{\alpha}^2(k-1).$$



Example : Primary colors are red, blue, yellow.

Every pick which is your favorite.

Identify

Red \Leftrightarrow 1

Blue \Leftrightarrow 2

Yellow \Leftrightarrow 3.

Random sample from class :

$$Y_1 = \underline{\hspace{2cm}}$$

$$Y_2 = \underline{\hspace{2cm}}$$

$$Y_3 = \underline{\hspace{2cm}}$$

Let's test

$$H_0 : P_1 = P_2 = P_3 = \frac{1}{3} .$$

$$\text{at } \alpha = 0.05.$$

$n = \underline{\hspace{2cm}}$ (# of people in class)

$$\chi^2_{\text{obs}} = \frac{(Y_1 - nP_{10})^2}{nP_{10}} + \frac{(Y_2 - nP_{20})^2}{nP_{20}} + \frac{(Y_3 - nP_{30})^2}{nP_{30}}$$

$$= + +$$

$$= + +$$

=

From the table,

$$\chi^2_{0.05}(2) = 5.991.$$

Therefore, we the null hypothesis.

The kind of χ^2 test we studied is sometimes called a " χ^2 goodness of fit" test.

There are other useful tests based on a χ^2 statistic.

Tests for Homogeneity

Let $X^1 =$ primary color chosen by a female
 $X^2 =$ " " " " " male.

Denote $P_{i1} = P(X^1 = i), i=1,2,3$

$P_{i2} = P(X^2 = i), i=1,2,3.$

Lets consider

$H_0 : P_{11} = P_{12}$ and $P_{21} = P_{22}$ and $P_{31} = P_{32}$

"Do males' preferences tend to differ from those of females?"

More generally, consider two experiments E_1 and E_2 , each with the same set of possible outcomes, $1, \dots, k$.

Define $P_{ij} = P(\text{experiment } j \text{ results in outcome } i)$

We want to test for

$$H_0 : P_{i1} = P_{iz} \text{ for all } i=1, \dots, k$$

Assume we have 2 random samples

$X_1^1, X_2^1, \dots, X_{n_1}^1$ from E_1

$X_1^2, X_2^2, \dots, X_{n_2}^2$ from E_2 .

Define $Y_{i1} = \# \text{ of times } i \text{ occurs in sample 1}$

$Y_{i2} = \# \text{ " " " " " " " 2.}$

Then

$$\sum_{i=1}^k \frac{(Y_{i1} - n_1 p_{i1})^2}{n_1 p_{i1}} \approx \chi^2(k-1)$$

$$\sum_{i=1}^k \frac{(Y_{i2} - n_2 p_{i2})^2}{n_2 p_{i2}} \approx \chi^2(k-1)$$

and so

$$\sum_{j=1}^2 \sum_{i=1}^k \frac{(Y_{ij} - n_j p_{ij})^2}{n_j p_{ij}} \approx \chi^2(2k-2)$$

Under H_0 , $p_{ii} = p_{i2}$ for each i .

In applications, however, these frequencies are unknown.

Thus it is customary to estimate $P_{i1} = P_{i2}$

with $\frac{Y_{i1} + Y_{i2}}{n_1 + n_2}$.

It can be shown that under H_0 ,

$$Q = \sum_{j=1}^2 \sum_{l=1}^k \frac{\left(Y_{ij} - n_j \cdot \frac{(Y_{i1} + Y_{i2})}{n_1 + n_2} \right)^2}{n_j \cdot \left(\frac{Y_{i1} + Y_{i2}}{n_1 + n_2} \right)}$$

$$\approx \chi^2(k-1)$$

($k-1$ degrees of freedom are lost in estimating the outcome frequencies).

Thus

$$Q \geq \chi^2_\alpha(k-1)$$

Reject H_0

has level of significance approximately α .

Back to our example (preferences for primary colors)

E1 : ask a female her favorite primary color

E2 : " " male his " " "

We observed

$$n_1 = \underline{\quad}$$

$$n_2 = \underline{\quad}$$

$$y_{11} = \underline{\quad}$$

$$y_{21} = \underline{\quad}$$

$$y_{12} = \underline{\quad}$$

$$y_{22} = \underline{\quad}$$

$$y_{13} = \underline{\quad}$$

$$y_{23} = \underline{\quad}$$

$$\hat{P}_1 = \frac{y_{11} + y_{12}}{n_1 + n_2} =$$

$$\hat{P}_2 = \frac{y_{21} + y_{22}}{n_1 + n_2} =$$

$$\hat{P}_3 = \frac{y_{31} + y_{32}}{n_1 + n_2} =$$

Then $Q = \sum_{j=1}^2 \sum_{i=1}^3 \frac{(Y_{ij} - n_i \hat{P}_i)^2}{n_i \hat{P}_i}$

$$= \frac{(Y_{11} - n_1 \hat{P}_1)^2}{n_1 \hat{P}_1} + \frac{(Y_{21} - n_2 \hat{P}_2)^2}{n_2 \hat{P}_2} + \frac{(Y_{31} - n_1 \hat{P}_3)^2}{n_1 \hat{P}_3}$$

$$+ \frac{(Y_{12} - n_2 \hat{P}_1)^2}{n_2 \hat{P}_1} + \frac{(Y_{22} - n_2 \hat{P}_2)^2}{n_2 \hat{P}_2} + \frac{(Y_{32} - n_2 \hat{P}_3)^2}{n_2 \hat{P}_3}$$

$$= \underline{\hspace{2cm}} + \underline{\hspace{2cm}} + \underline{\hspace{2cm}}$$

$$+ \underline{\hspace{2cm}} + \underline{\hspace{2cm}} + \underline{\hspace{2cm}}$$

$$= \underline{\hspace{2cm}} \quad \chi^2_{0.05}(2) = 5.991$$

Thus, we H_0 .

Contingency Tables

A contingency table is a concise way of summarizing data for testing homogeneity:

Suppose 500 people were polled on their opinion of smoking in public places:

Sex	Favor	Oppose	Indifferent	Totals
Male	151	73	12	236
Female	101	136	27	264
Totals	252	209	39	500

We can use this table to apply a χ^2 test for homogeneity.

Sometimes an investigator will ask

"are the responses /outcomes dependent on
the group?"

For example, "is preference for smoking in public
places dependent on sex?"

It turns out that testing

$$H_0 : P_i \cdot P_j = P_{ij} \quad (\text{independence})$$

results in the same test as testing for homogeneity.