

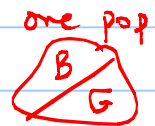
Lecture 22 (Ch. 8)

In Ch. 7, we learned how to build CIs for either 1 prop, π , or the difference between 2 props, $\pi_1 - \pi_2$, where π_1 = prop of something (e.g. boys) in population 1, and π_2 = " " same thing " " 2.

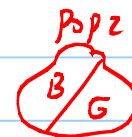
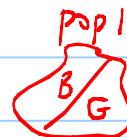
We also learned how to do hyp. tests on π , or $\pi_1 - \pi_2$.

[Note $\pi_1 + \pi_2 \neq 1$, because π_1, π_2 are 2 different populations]

But in all of these situations, the 2 pops have 2 categories (boy/girl) and π_i is the prop. of 1 of them.



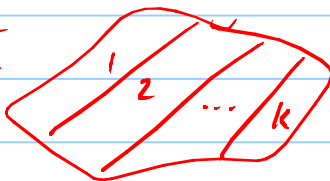
$$\pi_B = ?$$



$$\pi_{1B} - \pi_{2B} = ?$$

The tornado/climate eg. in prev. lecture deals with the situation where ONE population has 3 categories. For k categories:

one pop.



$$\pi_1 = ? , \pi_2 = ? , \dots , \pi_k = ?$$

We learned that the relevant dist. is chi-squared with $df = k - 1$. And the quantity that follows that dist. is

like z, t
$$\chi^2 = \sum_{i=1}^k \left(\frac{\text{obs}_i - \text{exp}_i}{\sqrt{\text{exp}_i}} \right)^2$$
 Order matters in interpretation (below).

where obs_i and exp_i are observed and expected counts in the i th category (still of 1 population). The latter are computed assuming H_0 is true, where

$$H_0 : \pi_1 = \pi_{01}, \pi_2 = \pi_{02}, \dots, \pi_k = \pi_{0k}$$

H_1 : At least one of these is wrong.

$$\left[\begin{array}{l} \text{This time } \pi_1 + \pi_2 + \dots = 1 \\ \pi_{01} + \pi_{02} + \dots = 1 \end{array} \right]$$

Note that the above H_0, H_1 is just a generalization of

$$H_0: \pi = \pi_0 \quad (z\text{-test}).$$

$$H_1: \pi \neq \pi_0$$

to more than 2 categories in the population.

However, there are no 1-sided / 2-sided varieties of chi-sq.

When X_{obs}^2 is small (say 10), then the observed counts are consistent with the expected counts if $H_0 = T$

(i.e. $\pi_1 = \pi_{01}, \pi_2 = \pi_{02}, \dots, \pi_k = \pi_{0k}$). So, if X_{obs}^2 is large,

then at least one of these must be wrong.

In other words the appropriate hypotheses are

$$H_0: \pi_1 = \pi_{01}, \pi_2 = \pi_{02}, \dots, \pi_k = \pi_{0k}$$

H_1 : At least one of these specifications is wrong.

And it is the "At least" which gives us

$$p\text{-value} = \text{prob}(X^2 > X_{obs}^2) \quad (\text{Table VII})$$

I.e. We are always interested in the upper tail area only

Said differently for the chi-sq test of the above H_0/H_1 , the p-value is only the right area, because violation of each part of H_0 , increases X^2 .

Interpretation/Diagnosis

The magnitudes of the k terms in χ^2_{obs} are important in deciding which of the k proportions are most different from the expected props. (under H_0).

Example: In the tornado example

Suppose we had found p -value $\leq \alpha$, i.e. There is evidence that climate does affect tornadic activity. Then χ^2_{obs} would be big. But what makes it big? The 3 terms contributing to χ^2_{obs} are:

El Nino

La Nina

Normal

1.27 (Large)

0.009

(Small)

0.49

→ Then we could conclude that it is the El Nino years which differ most (in terms of tornadic activity) from what we would expect under H_0 (i.e. if climate had no effect on tornadic activity).

And we could say that tornadic activity in La Nina years is pretty close to what one would expect by chance.

How about the "direction" of the association? E.g.

Are there more tornadoes in El Nino years than in Normal years?

Look at the signs of the k term is χ^2 , BEFORE squaring:

El Nino

La Nina

Normal

$$\chi^2_{obs} = \frac{(+4.9)^2}{18.9} + \frac{(-0.5)^2}{27.5} + \frac{(-4.4)^2}{39.6}$$

Note order!

In this formula, we looked at (expected - obs)²

So in El Nino years: $\text{exp.} > \text{obs} \Rightarrow$ Less tornadic than expected
in La Nina years: $\text{exp} < \text{obs} \Rightarrow$ More tornadic .. "

Q1: An information retrieval system has 3 storage locations, and it is designed with the expectation that the long-run proportion of requests for the 3 locations is $\frac{1}{4}, \frac{1}{2}, \frac{1}{4}$, respectively.

According to observations, however, of the 12 requests made over some period of time, the number of requests from the 3 locations is 2, 6, 4, respectively. We want to know if the design expectations are inconsistent with observations.

So, we do a chi-squared test. The value of χ^2_{obs} is

A)
$$\frac{\left(\frac{1}{4} - \frac{2}{12}\right)^2}{\frac{2}{12}} + \frac{\left(\frac{1}{2} - \frac{6}{12}\right)^2}{\frac{6}{12}} + \frac{\left(\frac{1}{4} - \frac{4}{12}\right)^2}{\frac{4}{12}}$$

B)
$$\frac{\left(\frac{12}{4} - 2\right)^2}{12/4} + \frac{\left(\frac{12}{2} - 6\right)^2}{12/2} + \frac{\left(\frac{12}{4} - 4\right)^2}{12/4}$$

C)
$$\frac{\left(\frac{1}{4} - 2\right)^2}{1/4} + \frac{\left(\frac{1}{2} - 6\right)^2}{1/2} + \frac{\left(\frac{1}{4} - 4\right)^2}{1/4}$$

D)
$$\frac{\left(\frac{12}{4} - \frac{2}{12}\right)^2}{12/4} + \frac{\left(\frac{12}{2} - \frac{6}{12}\right)^2}{12/2} + \frac{\left(\frac{12}{4} - \frac{4}{12}\right)^2}{12/4}$$

Counts not props.

$$\frac{(\text{exp.} - \text{obs})^2}{\text{exp.}}$$

The chi-sq d. distr. shows up in 2 other situations.

2) k props across r populations:

H_0 : r pops are homogeneous w.r.t. k categories

H_1 : not \uparrow big concept!

~~1/2/.../k~~ pop. 1

~~1/2/.../k~~ pop. 2

\vdots

~~1/2/.../k~~ pop. r

Homogeneous means

pop 1	π_1	π_2	π_3	\dots	π_k	$\rightarrow \sum_{i=1}^k \pi_i = 1$
pop 2	π_1	π_2	π_3	\dots	π_k	$\rightarrow \sum_{i=1}^k \pi_i = 1$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	
pop r	π_1	π_2	π_3	\dots	π_k	$\rightarrow \sum_{i=1}^k \pi_i = 1$

Not homogeneous means that at least 1 of these is wrong.

E.g. for $k=2$, The above H_0/H_1

Translate to:

H_0 : $\pi_1 = \pi_2 = \dots = \pi_r$

H_1 : At least 2 π 's are different

	Categ. 1	2
pop 1	π_1	$1 - \pi_1$
2	π_2	$1 - \pi_2$
\vdots	\vdots	\vdots
r	π_r	$1 - \pi_r$

homogeneous

[Note: This is diff. from $\pi_1 = \pi_{01}, \pi_2 = \pi_{02} \dots$]

3) Independence of 2 categorical variables,
one with k levels, The other with r levels.

Same
test,
diff.
conclusion

H_0 : 2 categ. vars. are indep.

H_1 : \dots not indep.

In such problems, the data are shown as a **Contingency Table**:

observed counts (from sample/data.)

of cases in categ 1 AND in pop 1

	category	1	2	3	row marginals
pop. A		a	b	c	a+b+c
B		d	e	f	d+e+f
		a+d	b+e	c+f	n ← total count
		column marginals			

The good news is that one can test homogeneity with a chi-sq dist., but with $df = (k-1)(r-1)$.

of cols
of categories } →
of rows
of populations

I.e. compute $\chi^2_{obs} = \sum_{\text{all cells}} \left(\frac{\text{obs} - \text{exp}}{\sqrt{\text{exp}}} \right)^2$, and p-value = $pr(\chi^2 > \chi^2_{obs})$

counts, not props

Table VII.

The only question is what are the expected counts?

Expected counts: Assuming $H_0 = T$.

	row-marginal	column marginal	
	$\frac{(a+b+c)(a+d)}{n}$	$\frac{(a+b+c)(b+e)}{n}$...
	$\frac{(d+e+f)(a+d)}{n}$

Remember this result a "row \times col. marginals".

(It's not easy to see this)

E.g. Are boys & girls homogeneous w.r.t. belief in afterlife?

Here are the data in
The form of a
Contingency Table:

	Yes	Un.	No	
Boys	435	58	89	582
Girls	375	50	84	509
	810	108	173	1109

Expected:
Counts

$$\begin{pmatrix} \frac{(582)(810)}{1109} & \dots & \dots \\ \dots & \dots & \dots \end{pmatrix} \begin{matrix} 582 \\ 509 \end{matrix}$$

810 108 173

$$= \begin{pmatrix} 432.1 & 57.6 & 92.3 \\ 377.9 & 50.4 & 80.7 \end{pmatrix}$$

$$\chi^2_{\text{obs.}} = \frac{(435 - 432.1)^2}{432.1} + \frac{(58 - 57.6)^2}{57.6} + \dots$$

$$= .019 + .0028 + \boxed{0.118} + \dots$$

0.022 + .0032 + $\boxed{0.135}$

$$= 0.3$$

See interpretation/diagnosis below.

Signs of (obs - exp)

$$\begin{pmatrix} + & + & - \\ - & - & + \end{pmatrix}$$

"big" numbers, both in the "No" category.

$df = (2-1)(3-1) = 2 \Rightarrow p\text{-value} > 0.1$ (huge)

↖ Table VII

Cannot reject H_0 in favor of H_1 , at $\alpha = .01$ (or .05)

I.e. there is no evidence to think that Boys and Girls are not homogeneous w.r.t. their belief in after life.

"different"

Also There is no evidence that Gender and Belief in afterlife are not independent. (i.e. They "are" independent.)

Interpretation / Diagnosis

So, based on this data, we cannot say that there is a difference between boys & girls w.r.t. their belief in afterlife.

Mathematically, the reason is that χ^2_{obs} was too small.

But suppose, χ^2_{obs} had turned out to be huge. Then we could conclude that there is a difference between boys & girls in terms of their belief in afterlife. Then just as before, we can look at the relative size of the various terms in χ^2_{obs} to see which ones make the χ^2_{obs} big.

In this example, the big terms are 0.118, 0.135, which correspond to the "No" category.

In short, if the result had turned out to be statistically significant (ie. $\chi^2_{obs} = \text{huge}$, $p\text{-value} < \alpha$), then we could go further and say that the biggest difference between boys and girls (in terms of their belief in afterlife) is in the non-believer category.

The signs can be interpreted, too; but we'll skip it for now.

Summary: Chi-sq shows up in 3 situations:

I) 1 pop. (1 variable) with k categories.

$H_0: \pi_1 = \pi_{01}, \pi_2 = \pi_{02}, \dots, \pi_k = \pi_{0k}$

H_1 : At least one of π is wrong.

Because there is only 1 pop, must have $\sum_{i=1}^k \pi_{0i} = 1$.

II) r pops. each with k categories:

pop 1	π_1	π_2	π_3	\dots	π_k	$\rightarrow \sum_{i=1}^k \pi_i = 1$
pop 2	π_1	π_2	π_3	\dots	π_k	$\rightarrow \sum_{i=1}^k \pi_i = 1$
\vdots	π_1	π_2	π_3	\dots	π_k	
pop r	π_1	π_2	π_3	\dots	π_k	$\rightarrow \sum_{i=1}^k \pi_i = 1$

H_0 : The above π 's are equal as specified.

(ie. the r pops are homogeneous w.r.t. the k categories)

H_1 : For at least 1 of the k categories, the proportions are not equal for all pops.

(ie. the r pops are not homogeneous ...)

III) Are 2 categorical variables independent?

H_0 : they are independent.

H_1 : they are not.

hw-lect 22-1

By hand

Have you ever wondered whether soccer players suffer adverse effects from hitting "headers"? The authors of the article "No Evidence of Impaired Neurocognitive Performance in Collegiate Soccer Players" (The Amer. J. of Sports Medicine, 2002: 157-162) investigated this issue. The paper reported that 45 of the 91 soccer players in their sample had suffered concussion, 28 of 96 nonsoccer athletes had suffered concussion, and only 8 of 53 student controls had suffered concussion. Denote

π_1 = pop. proportion of concussions among soccer players,

π_2 = pop. proportion of concussions among non-soccer players,

π_3 = pop. proportion of concussions among control group.

Set up this problem as a test of homogeneity of three populations with respect to 2 categories. Specifically,

- State the hypotheses in terms of π_1 , π_2 , π_3 .
- Write the data in the form of a contingency table.
- Compute the expected counts.
- Compute the p-value (or specify a range for it).
- State the conclusion "in English."
- Diagnose the various terms appearing in χ^2 .

hw-lect 22-2

By R

The accompanying data resulted from an experiment in which seeds of five different types were planted and the number that germinated within 5 weeks of planting was observed for each seed type ("Nondestructive Optical Methods of Food Quality Evaluation," Food Science and Nutr., 1984: 232-279). Carry out a chi-squared test at level .01 to see whether the proportion of seeds that germinate in the specified period varies according to type of seed.

#

# Seed type:	1	2	3	4	5
# Germinated:	31	57	87	52	10
# Failed to germinate:	7	33	60	44	19

#

Specifically,

- #
- Does the statement of the problem require a test of homogeneity of 2 populations with respect to 5 categories, or vice versa?
 - Compute the χ^2 , the df, and the p-value corresponding to your answer in part a.
 - State your conclusion "in English," at significance level 0.05.
 - Diagnose the magnitude of the various terms in χ^2

hw-lect 22-3

Consider 1 pop. with 2 categories (say A, B), and let π_A, π_B denote the proportion of A's and B's in the pop. Note $\pi_A + \pi_B = 1$.

In a hw problem (about Bell computers) you see that the 1-sample 2-sided z-test of $\begin{cases} H_0: \pi_A = \pi_0 \\ H_1: \pi_A \neq \pi_0 \end{cases}$ gives the same p-value as the

chi-squared test of proportions in 1 pop.

$$\begin{cases} H_0: \pi_A = \pi_0, \pi_B = (1 - \pi_0) \\ H_1: \text{At least one of these is wrong} \end{cases}$$

This equivalence can be seen at the level of H_0 and H_1 , too; note that "At least one of these is wrong" translates to $\pi_A \neq \pi_0$, because $\pi_A + \pi_B = 1$; and if $\pi_A = \pi_0$ is wrong, then so is $\pi_B = 1 - \pi_A$.

Now, consider 2 populations, each with 2 categories. Denote the 2 pops with "1", "2", and the 2 categories as "A", "B".

So, π_{1A} = prop. of A's in 1st pop., π_{2A} = ..., etc.

- Write down the H_0/H_1 that test homogeneity of pops w.r.t. categories
- Show that the H_0/H_1 in part a) are equivalent to H_0, H_1 of the 2-sample 2-sided z-test.

This document was created with Win2PDF available at <http://www.win2pdf.com>.
The unregistered version of Win2PDF is for evaluation or non-commercial use only.
This page will not be added after purchasing Win2PDF.