

Lecture 25 (Ch. 8)

In the procedure we have learned, the last step involves comparing the p-value with α . That practice is (slowly) becoming "old style". More recently, one reports the p-value itself, because by itself it's useful - it reflects the evidence from data against H_0 .

But, α does have an important interpretation nevertheless. We know that it is the largest prob at which we are confident to reject H_0 in favor of H_1 . But there is more to it!

Suppose we are testing $H_0: \mu \leq \mu_0$ vs. $H_1: \mu > \mu_0$.

We assume $H_0 = \text{True}$ (ie. $\mu = \mu_0$), then compute a p-value.

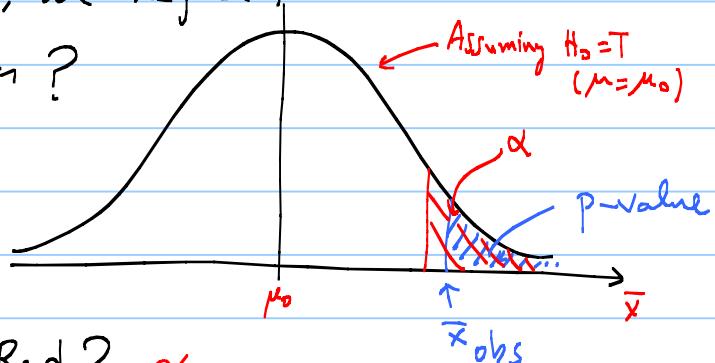
If p-value $< \alpha$, Then Reject H_0 in favor of H_1 .

So, every time p-value $< \alpha$, we reject.

How often will that happen?

For H_0, H_1 given here

$$\text{p-value} = \text{prob}(\bar{x} > \bar{x}_{\text{obs}})$$



Q How frequently is \bar{x} in the Red? α

Note: $\alpha = \text{prob}(\text{p-value} < \alpha \mid H_0 = T)$

So, $\alpha = \text{prob}(\text{Data Reject } H_0 \text{ in favor of } H_1 \mid H_0 = T)$

Type I error

"Bad" error

"False Alarm Rate"

(convicting an innocent person.)

This is how you decide on α .

How much bad error can you tolerate in the long run?

Why not set $\alpha = 0$, so that we will not have any bad (Type I) errors?

Because There is another kind of error :

$$\beta = \text{prob} (\text{Data cannot reject } H_0 \mid H_0 = \text{False})$$

in favor of H_1

Type II

(Releasing a guilty person.)

Setting $\alpha = 0 \Rightarrow \beta = 1$.

α, β (The probs of The 2 types of errors) have a complex but mostly inverse relationship, depending on n (p. 389-391)

So, given that α is The prob of The bad error, we generally set α at a fixed, but low, value.

Obviously, this will lead to some nonzero β , and it is important to compute it for your own specific problem.

Sometimes, people look at $1 - \beta$ (instead of β).

$$1 - \beta = \text{power} = \text{prob}(\text{Rejecting } H_0 \text{ --- } | H_0 = F)$$

(Convicting a guilty person).

If There is time, we'll return to power.

The understanding that $\alpha = \text{pr}(\text{Type I error}) = \text{pr}(\text{Bad error})$ offers another way of setting H_0/H_1 , correctly.

Example:

NASA

A company manufactures computer screens for used by astronauts on space missions. If more than 10% of the pixels on a given screen are defective, then the company does not release the screen to NASA, because otherwise disasters will occur.

Which is the appropriate H_0/H_1 ?

π = prop. of defective pixels. μ = mean # of defective pixels.

(A)

$$H_0: \pi < 0.1$$

(B)

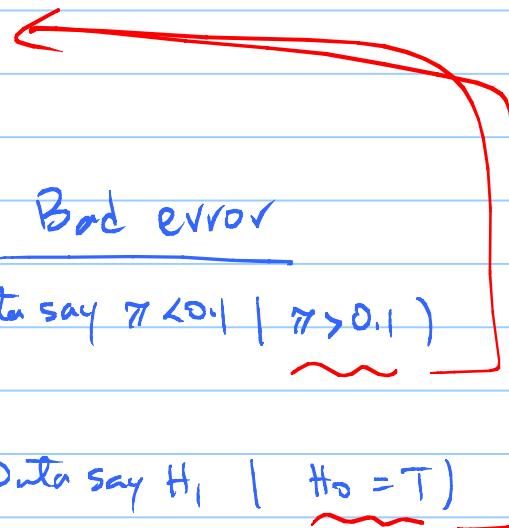
$$H_0: \pi = 0.1$$

$$H_1: \pi > 0.1$$

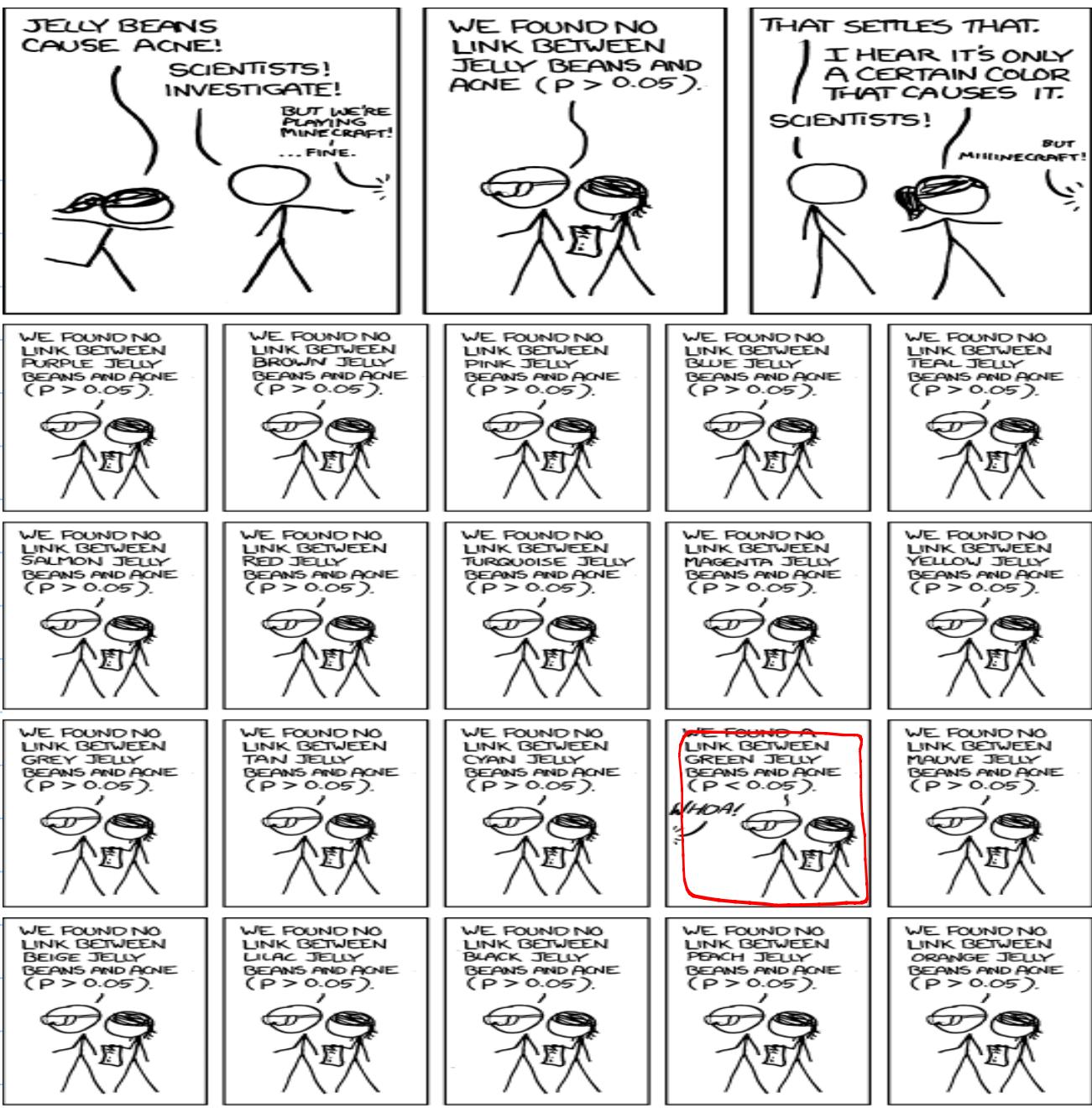
(C)

$$H_0: \pi > 0.1$$

$$H_1: \pi < 0.1$$



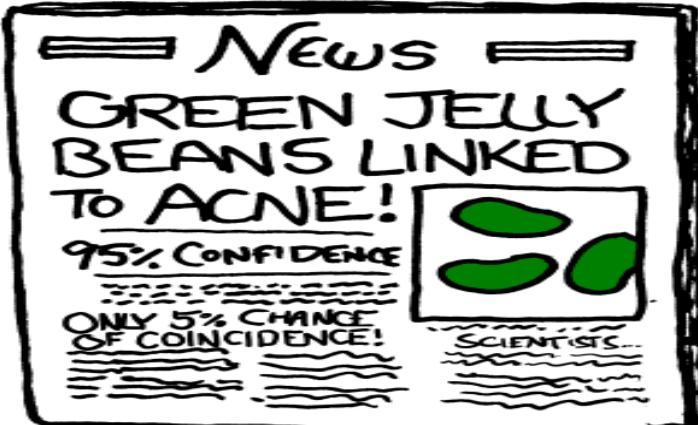
α
is
Dangerous!



Suppose you are testing whether a drug has $\mu > 0$.

So:

$$H_0: \mu \leq 0, H_1: \mu > 0$$



Suppose you compute the p-value and find $p\text{-value} > \alpha$, i.e. There is no evidence that $\mu > 0$. If you repeat the experiment many times, eventually you will find $p\text{-value} < \alpha$, i.e. There is evidence that $\mu > 0$. This will happen (at most) $\alpha\%$ of the time even if, in fact, $\mu < 0$.

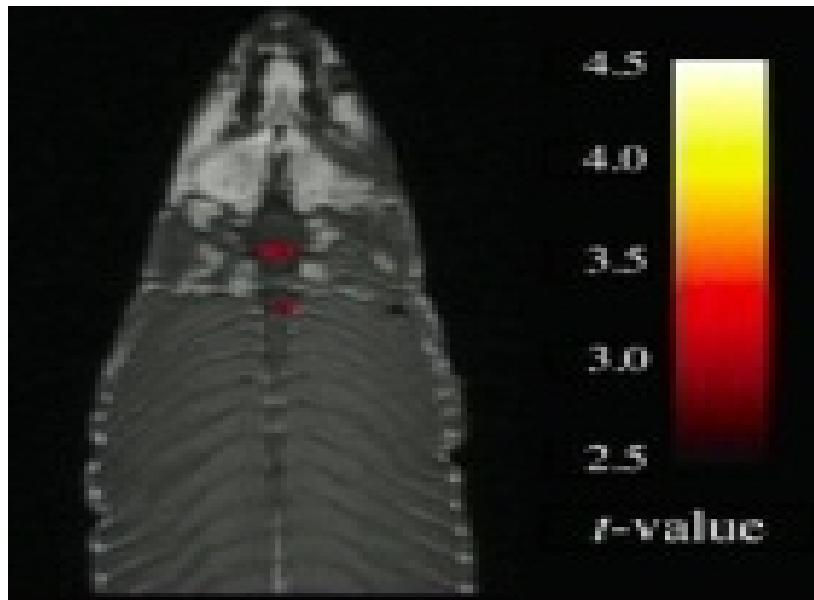
I.e. $\alpha\%$ of the time, you will make a type I error.

The conf. level (e.g. 0.95), and the significance level (e.g. 0.05) play an important role in every thing we've done. The last revealed what α really is : $\alpha = \text{prob}(\text{Type I}) = \text{prob}(\text{Bad error}) = \text{pr}(\text{Data rejects } H_0 \dots | H_0 = T)$

So, we now have another way of setting H_0 , H_1 , \rightarrow .

Also, we now see that we will be wrong (Bad wrong) $\alpha\%$ of the time. One example was the Jelly Bean example, here is another:

Dead Thinking Salmon!



There exist other decision-making frameworks which avoid such problems (e.g. check out

- multiple hypothesis testing
- False Discovery Rate)

Alternatively, in some situations, one can simply report the p-value, without comparing it to α .

In This class, we will continue to compare it with α , but be aware of this "defect"

FYI

Summary

We are done with 1-sample and 2-sample, z and t-tests, for paired and unpaired data, but all of that has dealt with the pop. means. What about pop. proportions?

Easy! Follow the pattern:

CI. for μ_x :

$$\bar{x} \pm z^* \frac{\sigma_x}{\sqrt{n}} \quad \text{or} \quad \left[\bar{x} \pm t^* \frac{s_x}{\sqrt{n}} \right] \quad df=n-1$$

C.I. for π_x :

$$p \pm z^* \sqrt{\frac{p(1-p)}{n}}$$

No p! ~~$p \pm z^* \sqrt{\frac{p(1-p)}{n}}$~~

Test for μ :

$$H_0: \mu = \mu_0 \quad H_1: \mu \neq \mu_0$$

$$z_{obs} = \frac{\bar{x}_{obs} - \mu_0}{\sigma_x / \sqrt{n}} \quad \text{or} \quad t_{obs} = \frac{\bar{x}_{obs} - \mu_0}{s_x / \sqrt{n}} \quad df=n-1$$

p-value = ...

Test for π :

$$H_0: \pi = \pi_0 \quad H_1: \pi \neq \pi_0$$

$$z_{obs} = \frac{p_{obs} - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}} \quad \text{Because we assumed } \pi = \pi_0.$$

p-value = ...

CI. for $\mu_2 - \mu_1$:

$$\bar{x}_2 - \bar{x}_1 \pm z^* \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad \left[\pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right] \quad df = \text{Welch}$$

Test for $\pi_2 - \pi_1$:

$$\pi_2 - \pi_1 \pm z^* \sqrt{\frac{\pi_1(1-\pi_1)}{n_1} + \frac{\pi_2(1-\pi_2)}{n_2}}$$

~~df = Welch~~

Test for $\mu_2 - \mu_1$:

$$H_0: \mu_2 - \mu_1 = \Delta \quad H_1: \dots$$

$$z_{obs} = \frac{(\bar{x}_2 - \bar{x}_1)_{obs} - \Delta}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad \text{or} \quad t_{obs} = \frac{(\bar{x}_2 - \bar{x}_1)_{obs} - \Delta}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

df = Welch

Test for $\pi_2 - \pi_1$:

$$H_0: \pi_2 - \pi_1 = \Delta \quad \dots$$

$$z_{obs} = \frac{(p_2 - p_1) - \Delta}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$$

When we compare 2 props (π_1, π_2), e.g. $H_0: \pi_1 - \pi_2 < 0.1$
 the implication is that we have TWO populations, each with
 2 groups/categ. (e.g. Boys, Girls).

In that case, ONE proportion (e.g. prop of boys, π_{Boys}) is enough
 to describe each pop., because the other prop. (e.g. π_{Girls}) is
 fixed by $1 - \pi_{\text{Boys}}$. The 2-sample z-test we have developed
 involves TWO proportions, one from each of TWO populations.

So, an example would be $\pi_1 = \pi_{\text{Boys}}$ in Northern hemisphere .

$$\pi_2 = \pi_{\text{Boys}} \text{ .. Southern ..}$$

Note that both π_1 and π_2 refer to Boys, but in 2 different
 populations (e.g. Northern and Southern hemispheres).

But there are situations where we have ONE population,
 with more than 2 categories, and we want to test
 some claim about the proportions of each category.

If we have ONE pop, with k categories, we can test

$$H_0: \pi_1 = \pi_{01}, \pi_2 = \pi_{02}, \dots, \pi_k = \pi_{0k} \quad \xrightarrow{\text{prop. of } k^{\text{th}} \text{ categ. in pop.}}$$

$H_1:$ At least one of π_i is wrong

I'll explain this later.

$$\sum_{i=1}^k \pi_{0i} = 1$$

Of course, given that there is only ONE pop., we have $\sum_{i=1}^k \pi_{0i} = 1$

Below, we will see how to do this test.

There will be a new distribution: Chi-squared.

Also note that a pop. with 2 groups can be thought of as being
 described by one random variable with 2 levels. Similarly, a pop.
 with k groups can be described with one r.v. with k levels.

E.g.

{Monthly Weather Review, 2008:
Vol. 136, p. 3121. Cook & Schaefer.

Does data provide sufficient evidence to support an association between climate and tornadic activity?

	El Nino	La Nina	Normal	
# of Days with violent tornadoes:	$n_1 = 14$	$n_2 = 28$	$n_3 = 44$	(86)
in each climate category				

$$\text{proportion: } \frac{14}{86} = 0.16 \quad 0.33 \quad 0.51 \quad (1)$$

Data.

# of years classified as	12	17	25	(54)
proportion:	$\frac{12}{54} = 0.22$	0.32	0.46	(1)

H_0 : true prop. of tornadic days in El Nino years. Etc.
There is no association, i.e.

$$H_0: \pi_1 = 0.22 \quad \pi_2 = 0.32 \quad \pi_3 = 0.46$$

H_1 : At least one of these assignments is wrong.

If H_0 = True, how many tornadoes do you expect in each of the $k=3$ categories?

$$\begin{array}{lll} \text{Expected} & 0.22(86) & 0.32(86) & 0.46(86) \\ \text{Count:} & \approx 18.9 & \approx 27.5 & \approx 39.6 \quad (86) \end{array}$$

$$\begin{array}{lll} \text{Observed} & 14 & 28 & 44 \\ \text{Count:} & & & \end{array}$$

$$(Exp - obs)^2: \quad (4.9)^2 \quad (-0.5)^2 \quad (-4.4)^2$$

$$\frac{(Exp - obs)^2}{Exp}: \quad 1.27 \quad 0.009 \quad 0.49$$

$$\text{Like } z_{obs}, t_{obs}, \chi^2_{obs.} = \sum_{i=1}^3 \frac{(exp - obs)^2}{exp} = 1.77$$

If there were really no difference at all in the # of tornadoes between the 3 categories, then this would be near zero.

Q So, is this χ^2_{obs} far away from 0 to reject H_0 (in favor of H_1)?
 Note: χ^2 is non-negative, unlike Z, t

A We need to know the samp. distr. of χ^2 , when $H_0 = T$.

Theorem: Under the null hypothesis, χ^2 has a chi-squared distr. with $df = k - 1$ ($= 3 - 1 = 2$)

What's a chi-squared dist? It's just another Table (VII).
 But FYI,

$$\text{p-value} = \text{prob}(\chi^2 > \chi^2_{\text{obs}}) = \text{prob}(\chi^2 > 1.77) > 0.1 \quad \begin{matrix} \uparrow \\ df = 3 - 1 = 2 \end{matrix} \quad \begin{matrix} \text{see a few} \\ \text{pages down} \end{matrix}$$

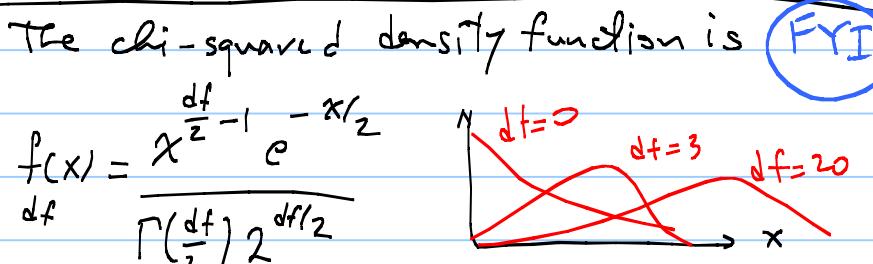
Conclusion (at $\alpha = .01$): p-value $> \alpha$ at least 1 is wrong.

In words: Cannot reject H_0 in favor of H_1 .
 $(\pi_1 = .22, \pi_2 = .32, \pi_3 = .46)$

In English: There is no evidence from data to suggest that the 3 props are NOT .22, .32, .46, i.e.

I.e. There is no evidence from data that there is an association between tornadic activity and climate.

For the chi-squared test, this sign is always $>$! See below for expl.



Summary / generalization

3, above

Now, let's generalize the above example to k categories:

Let π_i = proportion of cases in category i :

	Null params	Test results
π_1 = proportion of Categ. 1's	π_{01}	0.22
π_2 = - - -	π_{02}	0.32
π_3 = - - -	π_{03}	0.46

If $H_0 = \text{True}$, $H_0: \pi_1 = \pi_{01}, \pi_2 = \pi_{02}, \dots$

Then in a sample of size n , how many would

we expect in category 1 :	$n\pi_{01}$	18.9
" " " " 2 :	$n\pi_{02}$	27.5
" " " " 3 :	$n\pi_{03}$	39.6
⋮	$\sum_{i=1}^k n_i = n$	
But according to data, we observe this many :	n_1 n_2 n_3	14 28 44

Punch line:

How the theorem tells us that

$$\chi^2_{\text{obs}} = \sum_i \frac{(\text{exp.} - \text{obs})^2}{\text{exp.}} = \sum_{i=1}^k \frac{(n\pi_{0i} - n_i)^2}{n\pi_{0i}}$$

counts, not proportions!

has a chi-sqd. distr with $df = k-1$.

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

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

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

to more than 2 categories in the population.

) See how

However, there are [no] 1-sided / 2-sided varieties of chi-sqdf.

When χ^2_{obs} is small (say ≈ 0), then the observed counts are consistent with the expected counts if H_0 is true

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

then at least one of these ↑ must be wrong.

In other words the appropriate hypotheses are

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

$H_1:$ At least one of these ↑ specifications is wrong.

And it is the "At least" which gives us

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

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

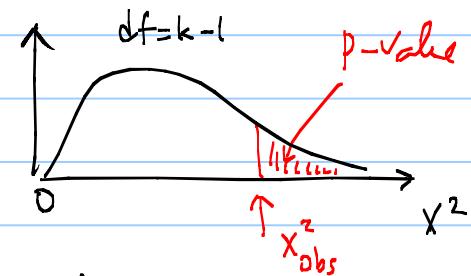
Said differently for the chi-sq test of the above H_0/H_1 , the p-value is only the right area, because violation of each part of H_0 , increases χ^2 .

How to use Table VII:

Table VII gives the area to the right of some value of χ^2_{obs} , i.e. it gives a p-value. However, it does not give all p-values; the only ones it provides are listed in the left-most column. E.g.

$$\chi^2_{\text{obs}} = 8.49, \text{ df}=4 \Rightarrow \text{p-value} = 0.075$$

$$\chi^2_{\text{obs}} = 8.66, \text{ df}=4 \Rightarrow \text{p-value} = 0.070$$



One might think that putting bounds on p-value is not enough for hypothesis testing, but it often is.

For example, suppose we get $\chi^2_{\text{obs}} = 8.55$ with $\text{df}=4$.

Then we can say $0.070 < \text{p-value} < 0.075$. That is good enough if $\alpha = 0.05$, because $\text{p-value} > \alpha$, and so we cannot reject H_0 in favor of H_1 .

FYI The same chi-sq test that we have developed above can be used to see if the row-variable and the column-variable in a matrix of counts are independent.

E.g. $\begin{array}{c|cc} & \overline{0} & \overline{Y} \\ \hline X | & \begin{array}{cc} 10 & 20 \\ 15 & 35 \end{array} \end{array}$ where 10 is the number cases with $x=0, Y=0$, etc.

This chi-squared test is called the test of homogeneity. Ask me, if you need to know more.

Summary

In Ch. 7, we learned how to build CIs for either 1 prop, π_1 , 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_x , 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.

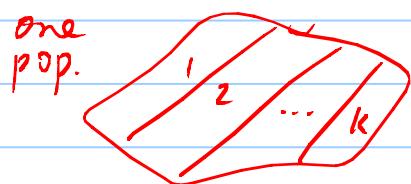


$$\pi_B = ?$$



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

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



$$\overline{r}_1 = ? , \overline{r}_2 = ? , \dots , \overline{r}_n = ?$$

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

$$X^2 = \sum_{i=1}^k \left(\frac{\text{obs}_i - \text{exp}_i}{\sqrt{\text{exp}_i}} \right)^2$$

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

$$H_0 : \pi_1 = \pi_{10}, \pi_2 = \pi_{20}, \dots, \pi_k = \pi_{k0} \quad \left[\text{This time } \pi_1 + \pi_2 + \dots = 1 \right]$$

H_1 : At least one of These^T is wrong.

$$\overline{D}_{01} + \overline{D}_{02} + \dots = 1$$

FYI

Further Diagnosis:

The magnitude of the k terms in χ^2_{obs} is useful for diagnosing which of the k proportions differ most from the Null values.

Example: In the tornado example

Suppose we had found p-value $<\alpha$, ie. There is evidence that climate does affect tornadic activity. In that case χ^2_{obs} would have been very large (so that p-value would be small).

The 3 terms contributing to χ^2_{obs} are

1.27	0.009	0.49
El Nino large	La Nina small	Normal

Then we could conclude that it is the El Nino years which are most different (in terms of tornadic activity) from what would be expected by chance (ie. if climate had no effect on tornadic activity).

In other words, we could conclude that the effect of climate on tornadoes is most in the El Nino years.

Tornadic activity in La Nina years, in fact, seems to be pretty close to what one would expect by chance.

Note: none of this tells us anything about the "direction" of the association. Are there more tornadoes in El Nino years than in Normal years? That's a different question that can be addressed by looking at the data directly. For example, one can look at # of tornadic days per El Nino years, per La Nina years, ... ($\frac{14}{12}, \frac{28}{17}, \frac{44}{25}$) = (1.17, 1.64, 1.76). El Nino years have "fewer" tornadoes.

hw-lect25-1

~~hw-lect24-1~~

hw-lect21-1 asked does it appear that π_x (The true proportion of defective screws) is at most 2.5%?

Then, the appropriate interval is the upper conf. Bound for π_x .

a) which of the following is/are appropriate pair of hypotheses.

A)

B)

C)

$$H_0: \pi_x \leq 2.5\%$$

$$H_0: \pi_x \geq 2.5\%$$

$$H_0: \pi_x = 2.5\%$$

$$H_1: \pi_x > 2.5\%$$

$$H_1: \pi_x < 2.5\%$$

$$H_1: \pi_x < 2.5\%$$

b) Compute the p-value (using the data in hw-lect21-1)

c) At $\alpha=0.05$, is the conclusion consistent with the conclusion from the CI approach in hw-lect21-1?

hw-lect25-2

A sample of 210 Bell computers has 56 defectives. Theory suggests that a third of all Bell computers should be defective. Does this data contradict the theory (at alpha=0.05)? Specifically,

a) Do a z-test ,

b) Do a chi-squared test with k=2 categories. Hint: The pi's (and pi_0's) of the k categories must sum to 1.

c) Are the conclusions in a and b consistent?

hw-lect25-3

Consider the data from an example in a past lecture where a survey of students in 390 yielded the following data:

17 students like Lab

48 " Do not like Lab

15 " have no opinion.

Suppose I believed that the proportion of students in each of the 3 categories (like, no-like, no-opinion) was equal.

Does this data contradict that belief? Let $\alpha=.05$.

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.