

Inference for Difference Between 2 Population Proportions

Chapter 10

We are interested in whether 2 pop. means are different.

Example:

$$\begin{aligned} \text{sample means} \\ \rightarrow \text{Mean Income in Bradford} &= \$38,010 \\ \rightarrow \text{Mean Income in Kane} &= \$35,006 \\ \text{difference} &= \$3,004 \end{aligned}$$

Chapter 11

We are interested in whether 2 pop. proportions are different.

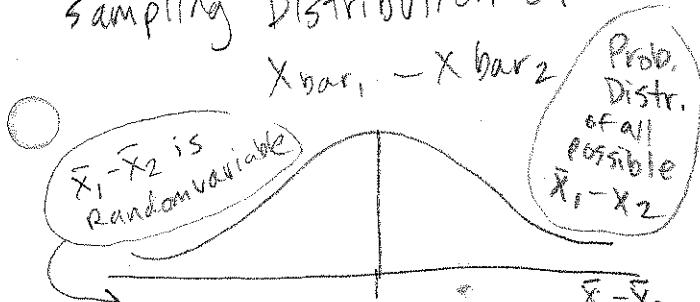
Example:

$$\begin{aligned} \text{sample proportions} \\ 2013 \text{ Data Entry Error Rate} &= 0.142 \\ 2014 \text{ Data Entry Error Rate} &= 0.107 \\ \text{difference} &= 0.035 \end{aligned}$$

Are the differences due to sampling error or are the significant differences that allow us to conclude that population parameters are different.

- ① Confidence Interval for difference between 2 Pop. proportions
- ② Hypothesis Test for difference between 2 pop. proportions
- ③ Hypothesis test to check Equality of 2 or more pop. proportions using chi-square Test statistic
- ④ Test of Independence of 2 categorical variables using chi-square test statistic

Chapter 10

Sampling Distribution of $\bar{X}_{\text{bar}_1} - \bar{X}_{\text{bar}_2}$ 

$$\begin{aligned}\text{mean} &= \mu_1 - \mu_2 \\ &= E(\bar{X}_{\text{bar}_1} - \bar{X}_{\text{bar}_2}) \\ &= \text{sum of all possible } \bar{X}_1 - \bar{X}_2\end{aligned}$$

Standard Error of Sampling Distribution of $\bar{X}_{\text{bar}_1} - \bar{X}_{\text{bar}_2}$

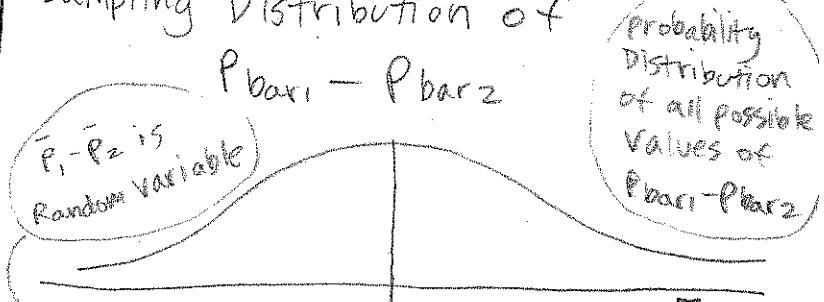
$$SD = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Correction Factor

$$\sqrt{\frac{N-n}{N-1}}$$

Needed when $\frac{n}{N} > 0.05$

Chapter 11

Sampling Distribution of $\bar{P}_{\text{bar}_1} - \bar{P}_{\text{bar}_2}$ 

$$\begin{aligned}\text{Mean} &= p_1 - p_2 \\ &= E(\bar{P}_{\text{bar}_1} - \bar{P}_{\text{bar}_2}) \\ &= \text{sum of all possible } \bar{P}_1 - \bar{P}_2\end{aligned}$$

Standard Deviation of Sampling Distribution of $\bar{P}_{\text{bar}_1} - \bar{P}_{\text{bar}_2}$

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

★ See Excel example to prove these formulas work.

In this chapter we won't know population proportions p_1 & p_2 so we will use \bar{p}_1 & \bar{p}_2 in formula to estimate $\bar{p}_1 - \bar{p}_2$

Remember Proportions from chapter 7:

$$\text{sample proportion} = \bar{p} = \frac{x}{n} = \text{Random variable}$$

P.2

x = number of elements in sample that possess the characteristic of interest

x = categorical variable

n = sample size

sampling distribution of \bar{p} can be approximated by a normal distribution whenever:

$$n * p \geq 5$$

$$n * (1-p) \geq 5$$

$$\left. \begin{array}{l} \text{Expected value of } \bar{p} \\ = E(\bar{p}) = p \end{array} \right\} \quad \begin{array}{l} \text{mean of all possible } \bar{p} \\ \text{Population proportion} \end{array}$$

Standard Error / Deviation of Sampling Distribution of \bar{p}

$$= \sigma_{\bar{p}} = \sqrt{\frac{p * (1-p)}{n}} * \sqrt{\frac{N-n}{N-1}}$$

Correction factor needed
when $\frac{n}{N} > 0.05$

① Confidence Interval for Difference Between 2 Population Proportions $P_1 - P_2$

P.3

- * From Sample Data, we calculate a point estimate for $P_1 - P_2$, $\bar{P}_1 - \bar{P}_2$, and then add a margin of error to both sides to get an interval (lower & upper value) that will contain the population proportion difference ($P_1 - P_2$) about 95 out of a 100 times.
- * If we don't know P_1 or P_2 , we use our sample proportions \bar{P}_1 and \bar{P}_2 in all of our calculations. (That is why we use word "Estimate" on SE calc.)
- * We must take 2 independent and random samples.
- * We can use the Normal Distribution to estimate the Sampling Distribution of $\bar{P}_1 - \bar{P}_2$ if 4 tests are passed:
 - ① $n_1 * \bar{P}_1 \geq 5$
 - ② $n_1 * (1 - \bar{P}_1) \geq 5$
 - ③ $n_2 * \bar{P}_2 \geq 5$
 - ④ $n_2 * (1 - \bar{P}_2) \geq 5$

Margin of Error

$$\underbrace{\bar{P}_1 - \bar{P}_2}_{\text{Point Estimate}} \pm \underbrace{Z_{\alpha/2} *}_{\# \text{ of Standard Deviations}}$$

$$\sqrt{\frac{\bar{P}_1 * (1 - \bar{P}_1)}{n_1} + \frac{\bar{P}_2 * (1 - \bar{P}_2)}{n_2}}$$

$P_1 - P_2$

of
 $P_1 - P_2$

$\# \text{ of Standard Deviations}$

Estimate of Standard Deviation/
Standard Error for Sampling
Distribution of $\bar{P}_1 - \bar{P}_2$

Variables Defined for Confidence Interval:

(P.4)

p_1 = population 1 proportion

p_2 = population 2 proportion

$\bar{p}_1 = p_{\text{bar}1}$ = Sample proportion from pop. 1

$\bar{p}_2 = p_{\text{bar}2}$ = Sample proportion from pop. 2

n_1 = sample size from population 1

n_2 = sample size from population 2

α = alpha = risk that the population proportion, $p_1 - p_2$, is not in interval

$1 - \alpha$ = confidence Level / Coefficient =
How sure we are that $p_1 - p_2$ is in interval

$\sqrt{\frac{N-n}{N-1}}$ = correction Factor for Standard Error calculation ie $\frac{n}{N} > 0.05$

(2)

Hypothesis Test to check if there is P.5
No Difference Between 2 population
proportions $P_1 - P_2$ ($\begin{matrix} \text{z method} \\ (\text{Hypothesized Difference} = 0) \end{matrix}$)

- * From Sample Data, we run our 5-step hypothesis Test "No Difference" between 2 population proportions.
- * In this section we run H.T. for $P_1 - P_2$ using the z method. Next section we will see how to check difference between 2 or more population proportions using Chi-square method. Chi-square can only do 2-tail Test.

* z method Allows us to run 3 types of Tests:

①	②	③
<u>1 tail, Lower</u>	<u>1 tail upper</u>	<u>2 tail</u>
$H_0: P_1 - P_2 \geq 0$	$H_0: P_1 - P_2 \leq 0$	$H_0: P_1 - P_2 = 0$
$H_a: P_1 - P_2 < 0$	$H_a: P_1 - P_2 > 0$	$H_a: P_1 - P_2 \neq 0$

* we must take 2 Independent & Random Samples

* We can use Normal Distribution to estimate the Sampling Distribution of $\bar{P}_1 - \bar{P}_2$ if 4 Tests are passed:

- ① $n_1 * \bar{P}_1 \geq 5$ ③ $n_2 * \bar{P}_2 \geq 5$
- ② $n_1 * (1 - \bar{P}_1) \geq 5$ ④ $n_2 * (1 - \bar{P}_2) \geq 5$

Formulas for Z method Hypothesis Test of $P_1 - P_2$

- ① IF H_0 is TRUE as an Equality:

$$P_1 - P_2 = 0$$

↓

$$P_1 = P_2$$

population proportion becomes: $P_1 = P_2 = P$

Estimate of

- ② Standard Error of $\bar{P}_1 - \bar{P}_2$ when $P_1 = P_2 = P$:

$$\sigma_{\bar{P}_1 - \bar{P}_2} = \sqrt{\frac{P_1 * (1-P_1)}{n_1} + \frac{P_2 * (1-P_2)}{n_2}} = \sqrt{P * (1-P) * \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

\downarrow
 $P_1 = P_2 = P \rightarrow \text{then}$ \uparrow

Now we just need estimate of P

- ③ Pooled Estimator of P when $P_1 = P_2 = P$

$\left\{ \begin{array}{l} \text{if } P_1 = P_2 = P \\ \text{we can just} \\ \text{take weighted} \\ \text{Average} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} \text{Pooled} \\ \text{Estimator} \\ \text{of } P \end{array} \right\} = \bar{P} = \frac{n_1 * \bar{P}_1 + n_2 * \bar{P}_2}{n_1 + n_2}$

weighted Average we learned in ch.3

- ④ Z Test statistic

$$Z = \frac{\bar{P}_1 - \bar{P}_2}{\sqrt{\bar{P} * (1-\bar{P}) * \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

point Estimate of Difference \leftarrow
SE \leftarrow

variables Defined for Hypothesis Test of $P_1 - P_2$

p_1 = population 1 proportion

(P.7)

p_2 = population 2 proportion

\bar{p}_1 = $p_{\text{bar}1}$ = sample proportion from Pop. 1

\bar{p}_2 = $p_{\text{bar}2}$ = sample proportion from Pop. 2

p = population proportion when $p_1 = p_2$

\bar{p} = estimate of p when we assume $p_1 = p_2 = p$,
called "Pooled Estimator of p " or weighted
Average of \bar{p}_1 and \bar{p}_2 .

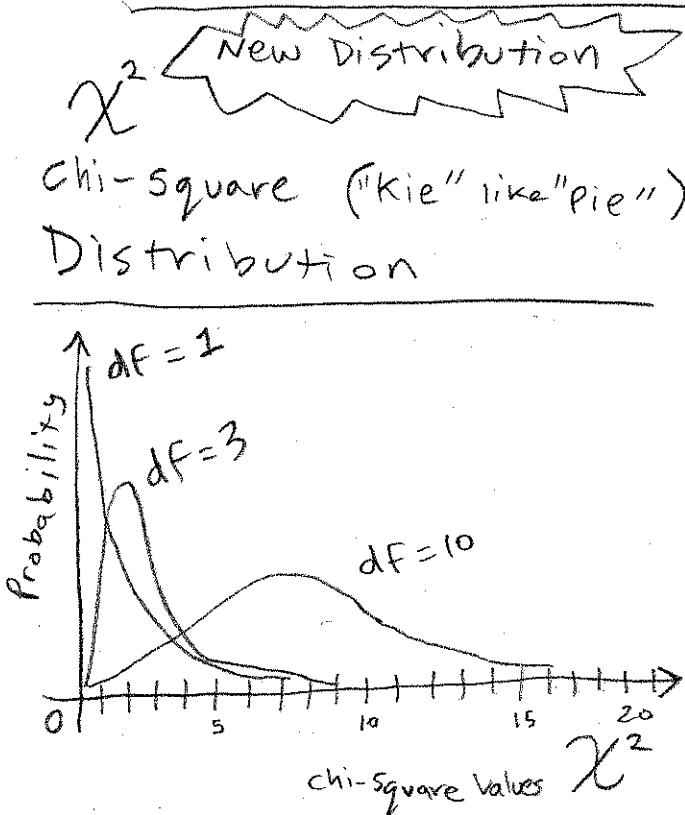
n_1 = sample size taken from population 1

n_2 = sample size taken from population 2

α = Risk that we Reject H_0 when it
is TRUE (Type I Error)

3 Hypothesis Test to check Equality of 2 or more population Proportions using Chi-Square

(P.10)



① Family of chi-square Distributions based on Degrees of Freedom, where:

K = # of populations

r = # of rows or categorical variables

For 2 categorical Variables	$df = K - 1$
For 3 or more categorical Variables	$df = (r - 1) * (K - 1)$

* called multinomial distribution
→ multiple populations

- ② Continuous Distribution where Area = Probability. All Area = 1
- ③ positive SKEW, but starts to approach Normal at $df = 10$
- ④ Because of calculations, χ^2 never negative. It is always an upper test.

χ^2 "Chi-Square" Test statistic

* using sample data, χ^2 test statistic can be used to test whether 2 or more population proportions are equal.

* χ^2 can be used to check Equality of 2 pop. proportions, but only for a 2-tail Test.

* To calculate the χ^2 test statistic and p-value for χ^2 we will compare "Observed Frequencies" for our proportions to the calculated "Expected Frequencies" for our proportions in a multiple step process that uses "Cross Tab" tables (Chapter 2) Remember: also known as "Contingency tables".

* The χ^2 test statistic calculations for "testing 2 or more pop. proportions" in this section is similar to the calculations for a "Test of Independence" in the next section. The difference is that "Test of Independence" has only 1 population, whereas "Testing 2 or more pop. proportions" has multiple populations.

* If you are checking Equality of 2 pop. proportions & you do Z method & χ^2 method, then: $(Z \text{ test statistic})^2 = \chi^2$

variables for χ^2 Hypothesis Test of Equality of 2 or More population proportions

P. 11

p_1 = population 1 proportion α = Alpha = Risk that we Reject H_0 when it is TRUE
 p_2 = population 2 proportion
 \vdots
 p_K = population K proportion

K = # of populations

r = # of rows in "Cross Tab" table or # of categorical variables.

\bar{p}_1 = p_{bar1} = Sample proportion from Pop. 1

\bar{p}_2 = p_{bar2} = Sample proportion from Pop. 2

\bar{p}_K = p_{barK} = Sample proportion from Pop. K

df = Degrees of Freedom

f_{ij} = observed frequency for Row i, column j in "Cross Tab" table

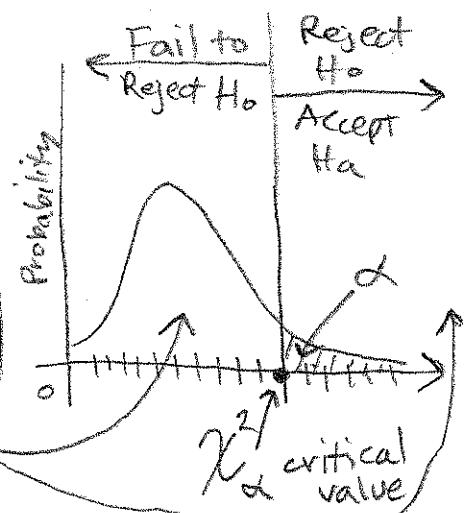
E_{ij} = Expected Frequency for Row i, column j

χ^2 = calculated Test statistic

χ^2_α = critical value (hurdle)

p-value = $P(\geq \chi^2)$

Is calculated
Test statistic
 χ^2 ?



Example of Hypothesis Test for Equality of 2 or more Population Proportions

P. 12

Question: Executive of HMO wanted to check the accuracy of patient contact information. A sample of patient records was taken for each of the years: 2011, 2012, 2013, 2014. Are error rates the same for each year?

Observed Frequencies

"Cross Tab" table with Frequencies

Sample Sizes

	2011 \bar{P}_1	2012 \bar{P}_2	2013 \bar{P}_3	2014 \bar{P}_4	overall P (yes)
Yes	$\frac{39}{343} = 0.114$	$\frac{43}{266} = 0.162$	$\frac{45}{316} = 0.142$	$\frac{41}{382} = 0.107$	$\frac{168}{1307} = 0.129$

- * If each \bar{P} (sample proportion) was equal to 0.129, then there would be zero difference (All Equal)
- * Remember, some difference (sampling error) is acceptable. we need to check for Significant Difference.

The essence of how we will check to see if there is a significant difference in the proportions is by comparing:

P.13

Observed Frequencies

and

Expected Frequencies

All observed Frequencies

Sample sizes

Errors / Years	2011	2012	2013	2014	Total
Yes	39	43	45	41	168
No	304	223	271	341	1139
Total	343	266	316	382	1307

* Remember, overall $P(\text{yes}) = \frac{168}{1307} = 0.129$

observed Frequency of "Yes" in 2013 = 45

sample size in 2013 = 316

$$2013 \bar{P}_3 = \frac{45}{316} = 0.142$$

* if $P(\text{yes in 2013})$ were to equal 0.129 then:

{sample size
2013}

$$* 0.129 = 316 * 0.129 =$$

40.6

check $\frac{40.6}{316} = 0.129$ ✓

"Expected Frequency"
if all proportions
equal!!!

5 Steps for Hypothesis Test to check Equality of 2 or more population proportions

P. 14

Step 1

Null & Alternative Hypothesis:

$$H_0: p_1 = p_2 = \dots = p_k$$

H_a : Not all population proportions are equal

1 or more pop. proportions differ from other pop. proportions

Step 2 a

Alpha

α = Risk that we reject H_0 when it was TRUE.

Step 2 b

Select Random Samples from each of the populations & create "cross Tab" table with Observed Frequencies.

Example:

"observed Frequencies" = f_{ij}

i = Row Number = Categorical variable number

j = Column Number = Population number

Pop. 1 Pop. 2 Pop. 3 Pop. 4

j=1 j=2 j=3 j=4

		Errors / Year 2	2011	2012	2013	2014	Total
		Yes	39	43	45	41	168
		No	304	223	271	341	1139
r=1	r=2	Total	343	266	316	382	1307

overall yes total

overall NOT total

Sum of all sample sizes

$f_{ij} = f_{13} = 45$ = observed Frequency Row 1, column 3

sample sizes

Step 3

Assume H_0 TRUE and compute
Expected Frequencies E_{ij}

P.15

E_{ij} = "Expected Frequencies"

i = Row Number = categorical variable Number

j = Column Number = Population Number

$$E_{ij} = \frac{\text{Row } i \text{ Total}}{\text{Sum of all sample sizes}} * \left(\frac{\text{Column } j}{\text{Total}} \right)$$

Observed Frequencies \rightarrow

Errors/Years	$j=1$	$j=2$	$j=3$	$j=4$	Row Totals
	2011	2012	2013	2014	Total
Yes	39	43	45	41	168
No	304	223	271	341	1139
Total	343	266	316	382	1307

$i=1$

$i=2$

Column Totals \rightarrow

$$E_{ij} = E_{13} = \frac{168}{1307} * 316 = 0.129 * 316 = 40.6$$

Grand over probability of getting "Yes"
 $P(\text{yes})$

Sample size for year
2013

expected value =
Expected count for "Yes"
in 2013 if proportion same
for each year

Completed Table of Expected Frequencies:

Expected Frequencies \rightarrow

Errors/Years	2011	2012	2013	2014	Total
Yes	44.1	34.2	40.6	49.1	168
No	298.9	231.8	275.4	332.9	1139
Total	343	266	316	382	1307

Step

4
acheck to see if all $e_{ij} \geq 5$ Expected Frequencies:TRUE All $e_{ij} \geq 5$

Errors/years	2011	2012	2013	2014
Yes	44.1	34.2	40.6	49.1
No	298.9	231.8	275.4	332.9

Step

4
bBecause All $e_{ij} \geq 5$, we can calculate:

{ Chi-Square Test Statistic }

$$\chi^2 = \sum_i \sum_j \frac{(f_{ij} - e_{ij})^2}{e_{ij}}$$

summation notation
for adding whole table with i rows &
j columns.

Example for Row i column j:

$$\frac{(f_{ij} - e_{ij})^2}{e_{ij}} = \frac{(f_{i3} - e_{i3})^2}{e_{i3}} = \frac{(45 - 40.62)^2}{40.62} = 0.4727$$

Full Table of calculation Answers for χ^2 :

Errors	years	2011	2012	2013	2014
Yes		0.587	2.269	0.4727	1.337
No		0.087	0.3347	0.0697	0.1971

Adding all these numbers = $\chi^2 = 5.35$

calculated test statistic

Step

4
C

Remaining Calculations:

P.17

○ $\text{Alpha} = \alpha = 0.05$

$K = \# \text{ of populations} = 4 \text{ years}$

$df = 4 - 1 = 3$ or if we use $(r-1)*(k-1)$
same answer: $(2-1)*(4-1) = 1*3 = 3$

$\chi^2_{\alpha} = \text{critical value} = \text{CHISQ.INV.RT}(\text{Alpha}, df)$
 $= \text{CHISQ.INV.RT}(0.05, 3) = 7.815$

$\left\{ \begin{array}{l} \text{P-value} \\ \text{Method 1} \end{array} \right\} = \text{CHISQ.TEST}(\text{actual_range, expected_range})$
 $= \text{CHISQ.TEST}(\text{Table of observed Frequencies, Table of Expected Frequencies})$
 $= 0.1476$

$\left\{ \begin{array}{l} \text{P-value} \\ \text{Method 2} \end{array} \right\} = \text{CHISQ.DIST.RT}(x, df)$

calculated χ^2 test statistic = 5.35 3

$= \text{CHISQ.DIST.RT}(5.35, 3) = 0.1476$

Step 5 Rules same as with previous Hypothesis Tests:

P-value $\leq \alpha$, then Reject H_0 & Accept H_a
critical value \geq calculated Test statistic, Reject.

Conclusion:
Statistical evidence does not suggest a significant difference between pop. proportions.

0.1476 \geq 0.05
Pvalue Alpha

5.35 \leq 7.815
calculated χ^2 critical v. χ^2

4) Test of Independence of 2 Categorical Variables using χ^2 chi-square ("Kie" like Pie) P.20

* Using Sample Data, we test if 2 categorical variables sampled from 1 population are independent (not dependent or associated).

* Example:

① Researchers want to determine if there is a difference in "Hiring Plans over the Next year" based on the "FIRM TYPE". Private Firm or Public Firm. (2 categorical variables)

② Human Resource Executives were surveyed and asked about their hiring plans over next year. The survey looked like this:

① what are your hiring plans over next year? Select one:

Hiring _____

Not Hiring _____

Lay off workers _____

one possible survey result might look like this

① what are your hiring plans over next year? Select one:

Hiring

Not Hiring _____

Lay off workers _____

② What sort of firm do you manage? Select one:

Private _____

Public

② What sort of firm do you manage? Select one:

Private _____

Public _____

③ From the random sample our survey results would look like:

Fields →
Records →

Hiring Plans over Next Year		FIRM TYPE
Hiring	Not Hiring	Public
Hiring		Private
Not Hiring		Public
Lay off workers		Private
:	:	:

2 categorical variables from 1 population of executives

P. 21

④ If objective of study is to determine if there is a difference in Hiring over Next Year Based on

when we cross tabulate to get Observed Frequencies

we must:

must put
"Response variable"
(Row variable)
"Hiring Plans over next year"
"Hiring Plans over next year"

MUST PUT
Explanatory Variable
("Column Variable")
"Firm Types"

		FIRM TYPE		Total
		Private	Public	
Plans over Next Year	Hiring	40	32	72
	Not Hiring	16	32	48
Hiring over Next Year	Lay off Employees	16	46	62
	Total	72	110	182

Sample size from 1 population

⑤ If we look at observed frequencies & proportions for just the "Response variable":

P.22

Hiring Plans over Next Year	observed Frequencies	overall sample proportion
Hiring	72	40%
Not Hiring	48	26%
Lay off Employees	62	34%
Total	182	100%

we can think of checking Independence of 2 categorical variables this way:

"Do the proportions for the "response variable" significantly differ when a second "explanatory variable" is added?"

2nd variable added

		Firm Type		
		Private	Public	Total
Hiring Plans over Next Year	Hiring	56%	29%	40%
	Not Hiring	22%	29%	26%
	Lay off workers	22%	42%	34%
	Total	100%	100%	100%

* is there a significant change?

* if there was zero change, they would all be the same as overall sample proportion.

⑥ ways we can ask Independence Question: p.23

- * Is the Response Variable "Hiring Plans over next year" independent of the Explanatory variable "Firm Type"?
- * Is the Response variable "Hiring Plans over Next Year" dependent on the Explanatory variable "Firm Type"?
- * Is there an association between the Response variable "Hiring Plans over Next Year" and the Explanatory Variable "Firm Type"?

⑦ our 2 possible conclusions:

① If 2 categorical variables are Independent:
* "Hiring Plans over Next Year" will not depend on "Firm Type" and the proportions for Hiring, Not Hiring and Laying off workers will be same for Public & Private Firms.

② If 2 categorical variables are NOT Independent:

- * Proportions for Hiring, Not Hiring & Lay off workers should differ by Firm.
- * After our chi-square test we will have evidence that "Hiring Plans over Next Year" variable is associated with or dependent on "Firm Type!"
- * We can gain more insights by comparing Sample proportions and making charts.

⑧ chi-square Test will involve same formulas as last section. we will compare "observed frequencies" to Expected Frequencies in a "cross tabulated table" sometimes called a "contingency Table":

* Assume H_0 TRUE & compute Expected Frequencies

$$e_{ij} = \frac{\text{Row } i \text{ Total}}{\text{Sample size}} * \frac{(\text{Column } j) \text{ total}}{\text{total}}$$

* If ALL $e_{ij} \geq 5$ then:

$$\left\{ \begin{array}{l} \text{chi-square} \\ \text{test} \\ \text{statistic} \end{array} \right\} = \chi^2 = \sum_{i=1}^n \sum_{j=1}^m \frac{(f_{ij} - e_{ij})^2}{e_{ij}}$$

④ Chi-square Test for Independence of 2 categorical variables (Contingency Table Test)

① Null & Alternative:

P.25

H_0 : The Response Variable (Row Variable) is independent of the Explanatory Variable (Column Variable).

H_a : The Response Variable (Row variable) is NOT independent of the Explanatory Variable (Column Variable).

② select alpha. Select Random sample from one population and collect data for both variables for every element in sample. Explanatory variable is the column variable & the response variable is the row variable for a cross tabulated Frequency table w/ observed frequencies = f_{ij} , where $i = \text{row}$ and $j = \text{column}$.

③ Assume H_0 TRUE & compute Expected Frequencies:

$$E_{ij} = \frac{\text{Row } i \text{ Total}}{\text{Sample Size}} * (\text{Column } j \text{ Total})$$

④ if all $E_{ij} \geq 5$, then compute chi-square Test statistic:

$$\chi^2 = \sum_{i,j} \frac{(f_{ij} - E_{ij})^2}{E_{ij}}$$

⑤ {Rejection
Rules}

Reject H_0 & Accept H_a if p-value $\leq \alpha$

Reject H_0 & Accept H_a if $\text{Test Statistic } \chi^2 \geq \text{critical value } \chi^2$

$$df = (r-1)*(c-1), \text{ where } r = \# \text{ Rows}$$

$$c = \# \text{ Columns}$$

⑥ If variables are NOT independent, investigate association by calculating Pbars & creating chart.

<p>Survey asked human resource executives about their firms hiring plans over the next year</p> <p>The survey offered three nominal "Hiring Over Next Year" categories: Hiring, Not hiring, Layoffs</p> <p>The survey also ask whether the firm was public or private. We will call this nominal variable:</p> <p>At the 0.05 alpha risk level, is the "Hiring Over Next Year" Nominal Variable independent of the "Firm Type" variable?</p> <p>Categorical Variables, or more specifically, Nominal Variables are:</p> <ul style="list-style-type: none"> "Firm Type" will be the "explanatory Variable" = "Column Variable" "Hiring Over Next Year" is the "Response Variable" = "Row Variable" <p>Different ways to ask Question about Independence:</p> <p>Is the "Hiring Over Next Year" variable independent of the "Firm Type" variable?</p> <p>Does the "Firm Type" variable influence/affect/dias the "Hiring Over Next Year" variable?</p> <p>Is the "Hiring Over Next Year" variable dependent on the "Firm Type" variable?</p> <p>Conduct a "Test of Independence" or "Contingency Table Test" to determine if the "Hiring Create "Cross Tab" or "Contingency Tables" for the "Observed Frequencies" and the "Expected Frequencies".</p>
0.05
H0 : The "Hiring Over Next Year" variable is independent of the "Firm Type" variable.
Ha : The "Hiring Over Next Year" variable is NOT independent of the "Firm Type" variable.
Alpha =
Chi-Square Test Statistic
Hiring Over Next Year/Firm Type
Hiring
Not Hiring
Lay Off Workers
Total
Private
Hiring
Not Hiring
Lay Off Workers
Total
Firm Type
Hiring
Not Hiring
Lay Off Workers
Total
Phi Proportions
Hiring Over Next Year/Firm Type
Hiring
Not Hiring
Lay Off Workers
Total
Expected Frequencies
Hiring Over Next Year/Firm Type
Hiring
Not Hiring
Lay Off Workers
Total
Chi-Square Test Statistic
Hiring Over Next Year/Firm Type
Hiring
Not Hiring
Lay Off Workers
Total
df = # Rows - 1 = 2
Rows = C = 3
Critical Value =
Chi-Square Critical Value =

<p>Categorical Variables, or more specifically, Nominal Variables are:</p> <p>"Firm Type" will be the "explanatory variable" = "Column Variable"</p> <p>"Hiring Over Next Year" is the "Response Variable" = "Row Variable"</p> <p>Different ways to ask Question about Independence:</p> <ul style="list-style-type: none"> Is the "Hiring Over Next Year" variable independent of the "Firm Type" variable? Does the "Firm Type" variable influence/affect/bias the "Hiring Over Next Year" variable? Is the "Hiring Over Next Year" variable dependent on the "Firm Type" variable? <p>Conduct a "test of independence" or "Contingency Table test" to determine if the "Hiring Over Next Year" variable is independent of the "Firm Type" variable.</p> <p>Create "Cross Tabs" or "Contingency Tables" for the "Observed Frequencies" and the "Expected Frequencies".</p>	<p>HO : The "Hiring Over Next Year" variable is independent of the "Firm Type" variable.</p> <p>Ha : The "Hiring Over Next Year" variable is NOT independent of the "Firm Type" variable.</p>																																																
<p>Alpha =</p>	<p>0.05</p>																																																
<p>Observed Frequencies</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="2">Hiring Over Next Year/Firm Type</th> <th>Total</th> </tr> <tr> <th>Hiring</th> <th>No Hiring</th> <th></th> </tr> </thead> <tbody> <tr> <td>Private</td> <td>Private</td> <td>=SUM(F23:F26)</td> </tr> <tr> <td>=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)</td> <td>=G23/G\$26</td> </tr> <tr> <td>Public</td> <td>Public</td> <td>=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$C\$1:\$C\$5)</td> </tr> <tr> <td>=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)</td> <td>=G24/G\$26</td> </tr> <tr> <td>Off Workers</td> <td>Off Workers</td> <td>=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$C\$1:\$C\$5)</td> </tr> <tr> <td>=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)</td> <td>=G25/G\$26</td> </tr> <tr> <td>Total</td> <td>Total</td> <td>=SUM(G23:G25)</td> </tr> </tbody> </table>	Hiring Over Next Year/Firm Type		Total	Hiring	No Hiring		Private	Private	=SUM(F23:F26)	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)	=G23/G\$26	Public	Public	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$C\$1:\$C\$5)	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)	=G24/G\$26	Off Workers	Off Workers	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$C\$1:\$C\$5)	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)	=G25/G\$26	Total	Total	=SUM(G23:G25)	<p>Expected Frequencies</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="2">Hiring Over Next Year/Firm Type</th> <th>Total</th> </tr> <tr> <th>Hiring</th> <th>No Hiring</th> <th></th> </tr> </thead> <tbody> <tr> <td>Private</td> <td>Private</td> <td>=F23/F\$26</td> </tr> <tr> <td>=SUM(F23:F26)</td> <td>=G23/G\$26</td> </tr> <tr> <td>Public</td> <td>Public</td> <td>=SUM(F24:F26)</td> </tr> <tr> <td>=SUM(F24:F26)</td> <td>=G24/G\$26</td> </tr> <tr> <td>Off Workers</td> <td>Off Workers</td> <td>=SUM(F25:F26)</td> </tr> <tr> <td>=SUM(F25:F26)</td> <td>=G25/G\$26</td> </tr> <tr> <td>Total</td> <td>Total</td> <td>=SUM(F23:F25)</td> </tr> </tbody> </table>	Hiring Over Next Year/Firm Type		Total	Hiring	No Hiring		Private	Private	=F23/F\$26	=SUM(F23:F26)	=G23/G\$26	Public	Public	=SUM(F24:F26)	=SUM(F24:F26)	=G24/G\$26	Off Workers	Off Workers	=SUM(F25:F26)	=SUM(F25:F26)	=G25/G\$26	Total	Total	=SUM(F23:F25)
Hiring Over Next Year/Firm Type		Total																																															
Hiring	No Hiring																																																
Private	Private	=SUM(F23:F26)																																															
=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)	=G23/G\$26																																																
Public	Public	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$C\$1:\$C\$5)																																															
=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)	=G24/G\$26																																																
Off Workers	Off Workers	=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$C\$1:\$C\$5)																																															
=COUNTIFS(\$A\$2:\$A\$163,\$E\$2:\$E\$5,\$B\$1:\$B\$5,\$D\$1:\$D\$5)	=G25/G\$26																																																
Total	Total	=SUM(G23:G25)																																															
Hiring Over Next Year/Firm Type		Total																																															
Hiring	No Hiring																																																
Private	Private	=F23/F\$26																																															
=SUM(F23:F26)	=G23/G\$26																																																
Public	Public	=SUM(F24:F26)																																															
=SUM(F24:F26)	=G24/G\$26																																																
Off Workers	Off Workers	=SUM(F25:F26)																																															
=SUM(F25:F26)	=G25/G\$26																																																
Total	Total	=SUM(F23:F25)																																															
<p>Chi-Square Test Statistic = $\chi^2 = \frac{\sum (O - E)^2}{E}$</p> <p>Chi-Square Critical Value =</p>	<p><== Must use Ctrl + Shift + Enter</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="2">Hiring Over Next Year/Firm Type</th> <th>Total</th> </tr> <tr> <th>Hiring</th> <th>No Hiring</th> <th></th> </tr> </thead> <tbody> <tr> <td>Private</td> <td>Private</td> <td>=SUM(F23:F26)</td> </tr> <tr> <td>=F23-F\$26</td> <td>=G23/G\$26</td> </tr> <tr> <td>Public</td> <td>Public</td> <td>=SUM(F24:F26)</td> </tr> <tr> <td>=F24-F\$26</td> <td>=G24/G\$26</td> </tr> <tr> <td>Off Workers</td> <td>Off Workers</td> <td>=SUM(F25:F26)</td> </tr> <tr> <td>=F25-F\$26</td> <td>=G25/G\$26</td> </tr> <tr> <td>Total</td> <td>Total</td> <td>=SUM(F23:F25)</td> </tr> </tbody> </table>	Hiring Over Next Year/Firm Type		Total	Hiring	No Hiring		Private	Private	=SUM(F23:F26)	=F23-F\$26	=G23/G\$26	Public	Public	=SUM(F24:F26)	=F24-F\$26	=G24/G\$26	Off Workers	Off Workers	=SUM(F25:F26)	=F25-F\$26	=G25/G\$26	Total	Total	=SUM(F23:F25)																								
Hiring Over Next Year/Firm Type		Total																																															
Hiring	No Hiring																																																
Private	Private	=SUM(F23:F26)																																															
=F23-F\$26	=G23/G\$26																																																
Public	Public	=SUM(F24:F26)																																															
=F24-F\$26	=G24/G\$26																																																
Off Workers	Off Workers	=SUM(F25:F26)																																															
=F25-F\$26	=G25/G\$26																																																
Total	Total	=SUM(F23:F25)																																															

Variables:

Turn Type

It looks like Private and Public

Add Employees is: 56% for

Lay Off Employees is: 23%

Hiring

56% 23%

Category	Percentage
Add Employees	56%
Lay Off Employees	23%

Turn Type

It looks like Private and Public

Add Employees is: 56% for

Lay Off Employees is: 23%

Hiring

56% 23%

Page 2 of 2 - Test of Independence

⑤ Chi-square Distribution & chi-square test statistic can be used in 3 situations:

P.28

Hypothesis Test for 2 or more population proportions

used to test whether 2 or more pop. proportions are equal.

$$\left\{ \begin{array}{l} \text{Expected Frequencies} \\ \text{under Assumption} \\ H_0 \text{ is TRUE} \end{array} \right\} = e_{ij} = \frac{\text{Row } i \text{ Total}}{\text{Sum of all sample sizes}} * (\text{Column } j \text{ Total})$$

$$\left\{ \begin{array}{l} \text{chi-square Test} \\ \text{Statistic if All} \\ e_{ij} \geq 5 \end{array} \right\} = \chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(f_{ij} - e_{ij})^2}{e_{ij}} \quad \begin{array}{l} \alpha = \text{alpha} \\ k = \# \text{populations} \\ c = \# \text{columns} \\ df = (r-1) * (k-1) \end{array}$$

3 or more row categorical variables (Multinomial Distr.) only 2 rows: df = k-1

$$\left\{ \begin{array}{l} \text{Critical Value for each} \\ \text{combination of 2 proportions} \\ - \text{done if test says there} \\ \text{is a difference} \end{array} \right\} CV_{ij} = \sqrt{\chi^2_{\alpha} * \sqrt{\frac{\bar{p}_i(1-\bar{p}_i)}{n_i}} + \sqrt{\frac{\bar{p}_j(1-\bar{p}_j)}{n_j}}}$$

Test of Independence "Contingency Table Test"

used to test if 2 categorical variables sampled from 1 population are independent. Does a set of proportions from the response variable (Row variable) significantly differ when a second explanatory variable (column variable) is added? $\rightarrow df = (r-1) * (c-1)$

$$e_{ij} = \frac{\text{Row } i \text{ Total}}{\text{sample size}} * (\text{Column } j \text{ Total})$$

Goodness of Fit Test (for nominal data)

used to determine whether a random ~~nominal~~ variable has a specific probability distribution. Check to see if observed frequencies for a multinomial probability distribution are equal to the expected frequencies.

$$\left\{ \begin{array}{l} \text{chi-square} \\ \text{test} \\ \text{statistic} \\ \text{if All} \\ e_i \geq 5 \end{array} \right\} = \chi^2 = \sum_{i=1}^K \frac{(f_i - e_i)^2}{e_i} \quad \begin{array}{l} f_i = \text{observed frequency} \\ \text{for category } i \\ e_i = \text{expected frequency} \\ \text{for category } i \\ K = \# \text{categories} \\ df = K-1 \end{array}$$

Independence Defined:

① Probability Theory (Chapter 4):

If occurrence of 1 event does not affect the probability of the other event.

② Probability Distributions (Chapter 5 on...)

2 Random Variables are Independent if the occurrence of one does not affect the probability Distribution of the other