The Difference in Proportions Test

Dr Tom Ilvento

Department of Food and Resource Economics



Decision Table for Two Proportions

- Use this table to help in the Difference of Proportions Test
- Large sample problems require that the combination of p*n or q*n > 5

Targets	Assumptions	Test Statistic
H ₀ : p ₁ - p ₂ = D	Independent Random Samples; Large sample sizes (n1 and n2 >50 when p or q > .10); Ho: D = 0	standard normal for comparisons with z; Pool the estimate of P based on the Null Hypothesis
	ndependent Random Samples; Large sample sizes (n1 and n2 >50 when p or q > .10); Ho: D≠ 0	standard normal for comparisons with z;

Overview

- A Difference of Proportions test is based on large sample only
- Same strategy as for the mean
 - We calculate the difference in the two sample proportions
 - Establish the sampling distribution for our estimator
 - Calculate a standard error of this sampling distribution
 - Conduct a test

2

Difference of two proportions

- For the null hypothesis
 - Ho: (P₁-P₂) = Do
 - Most often, Do = 0
- Since this is a large sample problem, we could use the sample estimates of p₁ and p₂ to estimate the standard error

$$S_{(p_1 - p_2)} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

 This is the standard error that we would use for a Confidence Interval for a difference of two proportions

%C.I. =
$$(p_1 - p_2) \pm z_{\alpha/2} s_{(p_1 - p_2)}$$

.

When Do = 0

- Under the Null Hypothesis where $P_1 P_2 = 0$
 - The Null doesn't specify what P is, so we don't know σ exactly
 - But the Null Hypothesis does tell us the proportions are equal, and thus the variances are equal
 - Thus it would be best to use information from both sample to estimate a single pooled variance
- We use a pooled average, P_p, based on adding the total number of successes and divide by the sum of the two sample sizes

$$P_p = \frac{(x_1 + x_2)}{(n_1 + n_2)}$$

• where x = # of successes for each group

5

When Do = 0

- In the case where the variances are assumed to be equal
 - We calculate the pooled proportion based on information from two samples
- $P_p = \frac{(x_1 + x_2)}{(n_1 + n_2)}$
- And then we use that pooled proportion in the calculation of the Standard Error

$$s_{(p_1-p_2)} = \sqrt{P_p Q_p \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

Note: this doesn't apply if Do ≠ 0, or when we calculate a Confidence Interval

6

Gender Differences in Shopping for Clothing

- A survey was taken of 500 people to determine various information concerning consumer behavior.
- One question asked, "Do you enjoy shopping for clothing?" The data are shown below by sex.
 - For women 224 out of 260 indicated Yes
 - For men, 136 out of 240 indicated yes
- Conduct a test that a higher proportion of women favor shopping using $\alpha = .01$

Start by establishing what you know from the problem

• $p_w = 224/260 = .8615$

• $p_m = 136/240 = .5667$

• Ho: $p_w - p_m = 0$ so we pool the variance

• $P_p = (224+136)/(260+240) = 360/500 = .72$

• $Q_p = (1-.72) = .28$

And the Standard Error?

• **SE** = $[.72*.28*(1/260+1/240)]^{.5}$ = .0402

7

Gender Differences in Shopping for Clothing

Ho:

Ha:

Assumptions

Test Statistic

Rejection Region

Calculation:

Conclusion:

• Ho: Pw - PM = 0

• Ha: Pw - Pm > 0 1-tailed test

· large samples; sigma known under Ho; pool

• $z^* = (.8615 - .5667 - 0)/.0402$

• $\alpha = .01$, z = 2.33

• $z^* = 7.335$

• z* > z

7.335 > 2.33

 Reject Ho: Pw - PM = 0 There is a difference!

.

pooled p 0.7200

П

Say it in words

- I conducted a test to determine if women indicated they enjoyed shopping for clothing more than men did.
- Because the Null Hypothesis was that the two groups are equal, I used a pooled estimate of the proportion who enjoyed shopping for clothing to calculate a pooled variance for the standard error of the test.
- For women, the proportion was .862 and for men it was .567
- The results of my test indicate a highly significant difference between women and men, p < .001. I have evidence to suggest that women enjoy shopping for clothing more than men.

Try a Difference of Proportions

Problem

• Geneticists have identified the E2F1 transcription factor as

10

The Excel file, DifMeans.xls

- You enter values wherever it is red
- It also does difference of means tests

Difference of Proportion Test Values in red need to be entered for the problem Group 1 136 Total 240 Ho: Difference 0.862 0.567 0.433 **Hypothesis Test** 0 2949 0 2949 Difference Difference Std Error Std Error 7.660 7.337 1-tailed 0.0000 1-tailed 0.0000 Confidence Interval 1.960 Z-value

0.219 to 0.370

0.075

an important component of cell proliferation control. The researchers induced DNA synthesis in two batches of serum-starved cells.
In one group of 92 cells (treatment), cells were microinjected with the E2F1 gene.
A control group of 158 cells was not exposed to E2F1.

A control group of 158 cells was not exposed to E2F1.
 After 30 hours, researchers determined the number of altered growth cells in each batch - use this to calculate a proportion.
 Test to see if the E2F1 treated cells had a higher proportion

of altered cells that the control group.

• Use α = .01

What calculations do you need to do?

	Altered	Not Altered	Row Total
E2F1	41	51	92
Control	15	143	158
Column Total	56	104	250

13

15

Altered Growth Cells

Ho:

Ha:

Assumptions

Test Statistic

Rejection Region

Calculation:

Conclusion:

• Ho: P_e - P_c = 0

• Ha: P_e - P_c > 0 1-tailed test

• large samples; sigma known under Ho; pool

• $z^* = (.446 - .095 - 0)/.0547$

• $\alpha = .01$, z = 2.33

• $z^* = 6.414$

• z* > z

6.414 > 2.33

• Reject Ho: Pe - Pc = 0

E2F1 has a higher proportion

Altered Growth Cells

Ho:

• Ha:

Assumptions

Test Statistic

• Rejection Region

Calculation:

Conclusion:

Difference of Proportion Test

Z-value

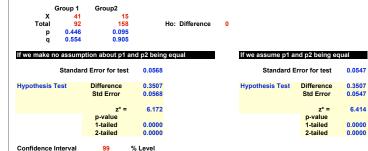
BOE

2.576 0.146 14

DifMeans.xls results

Values in red need to be entered for the problem

 Be sure you can find the information you need to conduct a test



0.204 to 0.497

16

pooled p 0.2240 pooled q 0.7760

Paired Difference Test

- When you have a situation where we record a pre and post test for the same individual, we cannot treat the samples as independent
- In these cases we can do a Paired Difference Test.
 - It's called "Matched Pairs Test" in some books.
 - Or **Mean Difference** in others
- The strategy is relatively simple
 - We create a new variable which is the difference of the pre-test from the post test
 - This new variable can be thought as a single random sample
 - For this new variable we calculate sample estimates of the mean and standard deviation

17

19

Paired Difference Test

- We calculate C.I. or conduct a hypothesis test on this new variable
- Often times the mean difference is referred to as
 - mean_D
- And the hypothesis is often:
 - Ho: $\mu_D = 0$
- This is no different than any single mean test, large sample or small sample

18

Sample data set for Paired Difference Test

Patient	Time I	Time 2	Difference tl - t2
1	5	5	0
2	I	3	
3	0	0	
4	I	ļ	
5	0	I	
6	2	I	

Alzheimer's study: Paired Difference Test

- Twenty Alzheimer patients were asked to spell 24 homophone pairs given in random order
 - Homophones are words that have the same pronunciation as another word with a different meaning and different spelling
 - Examples: nun and none; doe and dough
- · The number of confusions were recorded
- The test was repeated one year later
- The researchers posed the following question:
 - Do Alzheimer's patients show a significant increase in mean homophone confusion over time?
 - Use an alpha value of .05

20

Here are the sample statistics

- This can be thought of as a repeated measures study
- The sample size is 20

Statistics	Time 1	Time 2	Difference
Mean	4.15	5.80	1.65
Standard Error	0.78	0.94	0.72
Median	5.00	5.50	1.00
Mode	5.00	3.00	0.00
Standard Deviation	3.50	4.21	3.20
Sample Variance	12.24	17.75	10.24
Kurtosis	-0.85	0.13	0.08
Skewness	0.41	0.64	0.48
Range	11	16	12
Minimum	0	0	-3
Maximum	11	16	9
Sum	83	116	33
Count	20	20	20

21

Paired Difference Test

Ho:

Ha:

Assumptions

Test Statistic

Rejection Region

Calculation:

Conclusion:

22

Paired Difference Test

Ho:

Ha:

Assumptions

Test Statistic

Rejection Region

• Calculation:

Conclusion:

• Ho: $\mu_D = 0$

• Ha: $\mu_D > 0$ 1-tailed test

small sample; sigma unknown, use t

• $t^* = (1.65 - 0)/.72$

• $\alpha = .05$, $t_{.05, 19 df} = 1.729$

• t* = 2.29

• t* > t

• 2.29 > 1.729

• Reject Ho: μ_D = 0

• We find support for an increase 23

Difference of Means and Proportion tests

- First a difference or means or proportions is an extension of the basic hypothesis test
 - It may look more complicated
 - But it is the same thing!
- Begin by confirming it involves two random samples
- Then decide if it is two means or two proportions
- If two means can you assume equal variances?
- If proportions is the Null Hypothesis that the two groups are equal?
- We often rely on software to do these problems

24