

# Unit 6 - CI and hypothesis tests for proportions

## Confidence Intervals - Categorical

Confidence Interval means:

We are \_\_\_% confident that the true proportion of population is from \_\_\_ to \_\_\_.

It is the specified probability of a range including parameter.

### Conditions

- random
- Ind  $n < 0.1N$
- Normal  $n \cdot p \geq 10$   $n(1-p) \geq 10$

Main idea: We observe a sample. To guess the parameter, we make a range with the sample with degree of certainty because sample is likely close to the true proportion.

formula

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

margin error

Standard error

sample proportion

critical value

sample size

$z^*$ : Convert degree certainty into  $z^*$  on Calc:

Our C.I. rather catches pop parameter or not, with repeated samples+tests, our CI would capture it % certainty times.

- 2nd vars, Invnorm

- % is area, center cause CI, enter.

## Hypothesis tests - categorical

Main idea: We observe a sample. If the chances of getting that sample are significantly small with the null hypothesis being true, we reject the null.

### Conditions

- random
- $n \cdot p \geq 10$   $n(1-p) \geq 10$
- $n < 0.1N$

$\alpha$ -alpha

- Set alpha value: the point at which a smaller chance would be considered significant.

### Hypothesis

- Set hypotheses to know what rejecting null means.
- $H_0$ :  $H_0 = \text{Sign. previous claim}$
- $H_a$ :  $H_a: >, <, \text{ or } \neq$ , what we are inferring

C.I. difference in proportions

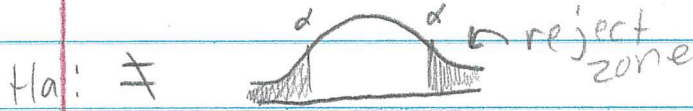
$$\hat{p}_1 - \hat{p}_2 \pm z^* \underbrace{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}_{\text{standard error Margin Error}}$$

If 0 is in the CI, it is not significantly different.

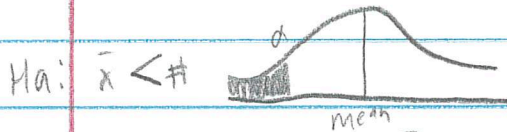
Solving P-value

formula: 
$$= \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = Z\text{-star}$$


plug into 2nd vars, norm cdf, draw picture to know which side to

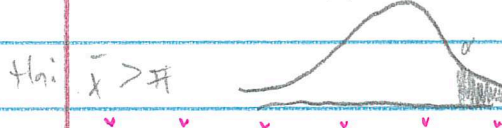


Split % ex) 95% → 97.25% put on.




P-Value Interpretation

1. When  $\alpha > p\text{-val}$   Significant reject null. When null is rejected,



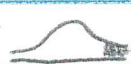
you can accept alternate

2. When  $\alpha < p\text{-val}$  

Logic: When null is rejected, the only other option is the alternate is true. When null is not rejected, we don't have enough information to accept it but we can't confirm that its false either.

Don't reject null, not significant, but don't accept null or alternate.

Error types

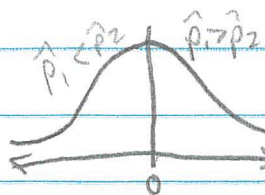
• E1: Reject a true null hypothesis probability:  $\alpha$   Ex)  $\alpha = .05$  means you are accepting a 5% chance you reject true

• E2: fail to reject false null hypothesis. probability: 1-power

• **power of test** is increased by

Differences hypothesis test

- Sample size increasing, or the
- $\alpha$  value increasing.



$$\frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}_c(1-\hat{p}_c) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$\hat{p}_c = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$$

→ Plug into 2nd vars, norm cdf, compare  $\alpha$  and p-value.

H<sub>0</sub>:  $\hat{p}_1 - \hat{p}_2 = 0$  no differ

H<sub>a</sub>:  $\hat{p}_1 - \hat{p}_2 <, >, \neq 0$

CI: If 0 fits in Interval, it is not significant diff



# Unit 7 - CIs and hypothesis tests for numerical data

## Conditions for all

random says

Normal  $n \geq 30$  or from normal pop

Independent  $n < .1N$

For Hypothesis tests  
always run Conditions,

alpha, Dof, and Hypothesis

For CI run conditions and dof.

## Numerical Confidence Intervals

Dof =  $n - 1$

$$t^* = \frac{\bar{x} - M}{\left(\frac{s}{\sqrt{n}}\right)}$$

$$CI = \bar{x} \pm t^* \left(\frac{s}{\sqrt{n}}\right)$$

or  $\text{InvT}_{\text{conf \%}, \text{dof}}$

Context: with \_\_\_% confidence the true parameter of \_\_\_ is from \_\_\_ to \_\_\_.

## Numerical Difference Confidence Intervals

Dof:  $(n_1 + n_2) - 2$

$$t^* = \frac{\bar{x} - M}{\left(\frac{s}{\sqrt{n}}\right)}$$

or  $\text{InvT}_{\text{dof}, \%}$

$$C.I. = (\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

With \_\_\_% confidence, the true difference in means of \_\_\_ and \_\_\_ is from \_\_\_ to \_\_\_.

## Numerical Hypothesis test

Do Conditions, Alpha, Hypothesis, and Dof.

Dof:  $n - 1$

$\alpha =$  usually  $.05$

$$t^* = \frac{\text{obs} - \text{exp}}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

is included in the interval, there is not a significant difference between the two. **OR**. Since 0 isn't in our interval there's a significant difference.

$\rightarrow$  tcdf  $\rightarrow$  p-val  $\rightarrow$  compare to reject/not reject

## Hypothesis test numerical Differences

Dof: run these two Separately, test both

smaller sample  $\rightarrow$

$n_1 - 1$  and  $(n_1 + n_2) - 2$

$$t^* = \frac{(\bar{x}_1 - \bar{x}_2) - (M_1 - M_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

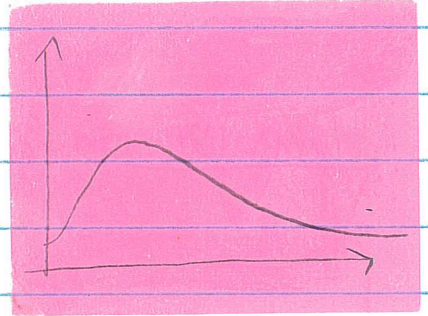
$\rightarrow$  tcdf, compare p-value

## Unit 8 - $\chi^2$ Homogeneity and Independence tests

We use  $\chi^2$  tests to observe if different variables are independent or dependent, and homogeneous/heterogeneous

Chi-Squared Dist:

- Almost Always right skewed
- Gets more normal as sample size goes up



Always do CHAD (conditions, hypothesis, alpha, dof)

### Conditions

- random says
- normal  $n \cdot p \geq 5$  for all expected values
- Independent  $n < 1N$

### Dof

1 Var:  $(\# \text{ columns} - 1)$   
 2 vars:  $(\text{rows} - 1) \cdot (\text{columns} - 1)$

$\alpha$  = usually .05

unless says otherwise

### Hypotheses

$H_0$ : <sup>homo</sup>no differ dist, <sup>independence</sup>no association  
 $H_a$ : differ dist, association

### Differences of Homogeneous and Independent

homogeneous -

- Compares 1 variable across 2 populations
- Asks if difference in the distribution of populations

Independent -

- Compares 2 variables for one population
- Asks if association between 2 variables in a population

### Main order of steps

1. Find expected values, plug in obs data L1, exp data L2
2. CHAD, plug lists into formula, Sto  $\rightarrow$  L3, Stats Calc vars on L3, Look at  $\chi^2$ . That's the test statistic  $\rightarrow$  2nd vars  $\chi^2$  cdf, <sup>compare</sup> p-value



## Calculating $\chi^2$

Expected values -

- If %'s, multiply each % by the total sample size.
- If table, do  $\frac{(\text{row total} \cdot \text{column total})}{\text{total}}$

Ex)

	kayak	bike	raft	
NOC	79 observed			720
WW				
Gorge				
	123			1745

To find the expected # of kayakers at NOC, do:

$$\left( \frac{720 \cdot 123}{1745} \right)$$

• Plug obs into  $L_1$ , expected into  $L_2$

$$\chi^2 = \sum \frac{(\text{obs} - \text{exp})^2}{\text{exp}} \quad \text{So...} \quad L_3 = \frac{(L_1 - L_2)^2}{L_2}$$

Sto that into  $L_3$ , then calc  $L_3$  and  $\sum x$  is test stat.

Convert to p-value!  $\sum x \rightarrow$  2nd vars, #8  $\chi^2$  cdf.

If  $p > \alpha$  not reject null, if  $p < \alpha$  reject null

results not significant; accept alt of associated  
no association between vars; or difference in dist. of  
or no difference in dist. of populations.

## Residuals - part of unit 9

• It is possible to make residual points from  $L_2 - L_3$  or (obs - exp)

To get  $L_3$ :  $y = \text{slope}(L_1) + y\text{-int}$ . This gives expected values of LSRL. This tests if the LSRL is a good fit. Do  $(L_1, L_4)$  as residual points. If no pattern, linear is good fit.

# Unit 9 - Slope applied to CI and Hypothesis tests

## 2-variables

How to explain slope in context:

The y-var: \_\_\_\_\_ goes up/down \_\_\_\_\_ units for every one x-var: \_\_\_\_\_ goes up.

There is variation between slope and y-int of different samples, so investigate with CI + hypo test.

## C.I with slope

Formula:

$$b \pm t^* (SE)$$

Dof =  $n-2$

$t^*$ : InvT, %, Dof

This C.I. gives range of true slope.

Context: we are \_\_\_\_\_ % conf  
The true population slope of \_\_\_\_\_ is from \_\_\_\_\_ to \_\_\_\_\_.

If 0 is in our interval, it could be reasonable to conclude there is no association between the two-variables.

## Computer outputs

variable	Coef	SE	t-ratio	prob
Constant	# y-int	<sup>standard error</sup> y-int		
variable	# slope	<sup>standard error</sup> slope		
	$r^2 = \#$	$S = \#$		

## Variables

$\beta$  = slope pop       $A$  = y-int pop  
 $b$  = slope sample       $a$  = y-int sample

Prob row: the probability we get that sample if true slope = 0.  
So low prob = high chance of association

# THE END!