

Unit 5: Inference for categorical data

1. Inference for a single proportion

Sta 101 – Spring 2019

Duke University, Department of Statistical Science



Dr. Ellison

Slides posted at
<https://www2.stat.duke.edu/courses/Spring19/sta101.001/index.html>

Outline

1. Housekeeping

2. Main ideas

Hypothesis Testing and Confidence Intervals with a Single Categorical Variable

1. Theoretical Hypothesis Testing and Confidence Intervals: \mathcal{Q} The CLT also describes the distribution of \hat{p}
2. Interpretations: \mathcal{M} \mathcal{Q} CI vs. HT determines observed vs. expected counts / proportions
3. CLT Conditions: \mathcal{Q} Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution

3. Applications

1. \mathcal{H} Single population proportion, large sample
2. \mathcal{H} \mathcal{M} Single population proportion, small sample

4. Recap

5. Summary

Organizing Analyses We Know

Tip: Break it down by types of variables involved in the research question.



<https://www.facebook.com/konmarimethod/>

Organizing Analyses We Know

Tip: Break it down by types of variables involved in the research question.

- Any **numerical** variables? How many?
- Any **categorical** variables? How many?

Organizing Analyses We Know

Tip: Break it down by types of variables involved in the research question.

- Any **numerical** variables? How many?
- Any **categorical** variables? How many?
 - How many **levels** does each categorical variable have?

Organizing Analyses We Know

Tip: Break it down by types of variables involved in the research question.

- Any **numerical** variables? How many?
- Any **categorical** variables? How many?
 - How many **levels** does each categorical variable have?
- If two variables, which variable(s) are **explanatory** and which variable is **response**?

Organizing

What analyses do we know so far?



Organizing

What analyses do we know so far?

Is there more than one way to conduct them?

Organizing

What analyses do we know so far?

Is there more than one way to conduct them?

If so, when should we use each way (ie. what conditions must be met)?

Organizing

Types of analyses that **HAVE one specified Population Parameter of Interest** and

Organizing

Types of analyses that **HAVE one specified Population Parameter of Interest** and

- We can create a **confidence interval** for the population parameter.

Organizing

Types of analyses that **HAVE one specified Population Parameter of Interest** and

- We can create a **confidence interval** for the population parameter.
- We can conduct a **hypothesis test** for the population parameter using:

Ho: Pop. Param = #

Ha: Pop. Param (\neq or $<$ or $>$) #

Organizing (read over later)


What's new in Unit 5?

Types of Variable(s) Involved	Population Parameter	Analyses and Ways to Conduct Them (that we know so far)	
		Confidence Interval for the Population Parameter	Hypothesis Test for the Population Parameter <i>H₀: Pop. Param = #</i> <i>H_a: Pop. Param (< or >) #</i>
Single Numerical Variable	μ	CLT Confidence Interval (Unit 3+4) Bootstrap Confidence Interval (Unit 4)	CLT Hypothesis Test (Unit 3+4) Bootstrap Hypothesis Test (Unit 4)
	μ_{diff}	CLT Confidence Interval (Unit 4) Bootstrap Confidence Interval (Unit 4)	CLT Hypothesis Test (Unit 4) Bootstrap Hypothesis Test (Unit 4)
	Median	Bootstrap Confidence Interval (Unit 4)	Bootstrap Hypothesis Test (Unit 4)
Single Categorical Variable (2 levels)	p	CLT Confidence Interval (Unit 5) Bootstrap Confidence Interval (Unit 4)	CLT Hypothesis Test (Unit 5) Bootstrap Hypothesis Test (Unit 4) Randomization Testing (Unit 5-Selecting balloons out of bag, rolling dice)
Numerical Response Variable Categorical Explanatory Variable (2 levels)	$\mu_1 - \mu_2$	CLT Confidence Interval (Unit 4) Bootstrap Confidence Interval (Unit 5)	CLT Hypothesis Test (Unit 4) Randomization Testing (Unit 1-Shuffling Cards)
	Median1-Median2	Bootstrap Confidence Interval (Unit 5)	Randomization Testing (Unit 1-Shuffling Cards)
Categorical Response Variable Categorical Explanatory Variable (both have 2 levels)	$p_1 - p_2$	CLT Confidence Interval (Unit 5) Bootstrap Confidence Interval (Unit 5)	CLT Hypothesis Test (Unit 5) Randomization Testing (Unit 1-Shuffling Cards)

Organizing

Types of analyses that **DON'T HAVE one specified Population Parameter of Interest** and

- No confidence interval framework
- Hypothesis tests set up are different



<https://www.facebook.com/yourmethod/>


Organizing (read over later)

What's new in Unit 5?


Types of Variables	Analysis	Hypotheses
Numerical Response Variable Categorical Explanatory Variable (>2 levels)	ANOVA	$H_0: \mu_1 = \mu_2 = \dots = \mu_k$ H_a : at least one pair of groups has means that are different
Single Categorical Variable (>2 levels)	Chi-Squared Goodness of Fit Test	H_0 : The data follows the specified distribution. H_a : The data does not follow the specified distribution.
Categorical Response Variable Categorical Explanatory Variable (at least one has >2 levels)	Chi-Squared Independence Test	H_0 : The two variables are independent/not associated. H_a : The two variables are dependent/associated.

Focus today:

Categorical Variable (2 levels)
Single Population Proportion, p
 $p = \text{proportion of one of the levels}$




Focus today:
 Categorical Variable (2 levels)
 Single Population Proportion, p
 $p = \text{proportion of one of the levels}$




- Construct [confidence intervals](#) and [hypothesis tests](#) using **Central Limit Theorem methods**.

Focus today:
 Categorical Variable (2 levels)
 Single Population Proportion, p
 $p = \text{proportion of one of the levels}$



- Construct [confidence intervals](#) and [hypothesis tests](#) using **Central Limit Theorem methods**.
- Conduct a [hypothesis test](#) for p with **randomization testing**.
- Create a [confidence interval](#) for p with **bootstrapping**.

Focus today:
 Categorical Variable (2 levels)
 Single Population Proportion, p
 $p = \text{proportion of one of the levels}$



When to use each method

- Construct [confidence intervals](#) and [hypothesis tests](#) using **Central Limit Theorem methods**.
- Conduct a [hypothesis test](#) for p with **randomization testing**.
- Create a [confidence interval](#) for p with **bootstrapping**.

Outline

- Housekeeping
- Main ideas
 - Hypothesis Testing and Confidence Intervals with a Single Categorical Variable**
 - Theoretical Hypothesis Testing and Confidence Intervals: \mathcal{Q} The CLT also describes the distribution of \hat{p}
 - Interpretations: new \otimes CI vs. HT determines observed vs. expected counts / proportions
 - CLT Conditions: \mathcal{Q} Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution
- Applications
 - new Single population proportion, large sample
 - new new new Single population proportion, small sample
- Recap
- Summary


🔍
Distribution of \hat{p}

**Organizing CLT Results
for Different Sample
Statistics**

When other **certain conditions** are met...
 $\bar{x} \sim N(\text{mean} = \mu, \text{standard dev./error} = \frac{\sigma}{\sqrt{n}})$

When other **certain conditions** are met...
 $\bar{x}_1 - \bar{x}_2 \sim N(\text{mean} = \mu_1 - \mu_2, \text{standard dev./error} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}})$

When other **certain conditions** are met...
 $\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$



https://www.facebook.com/konmarimethod/

🔍
Distribution of \hat{p}

Central Limit Theorem for Proportions
When **certain conditions** are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

Conditions:

- ▶ **Independence:**
 - a. Random sample/assignment AND
 - b. $n < 10\%$ of population
- ▶ **Shape/Skew:**
 - a. "Success/Failure Conditions:"
 - a. $np \geq 10$ AND
 - b. $n(1-p) \geq 10$

🔍
Distribution of \hat{p}

Why are the following called the **success/failure** conditions?

$$np \geq 10$$

$$n(1-p) \geq 10$$

🔍
Distribution of \hat{p}

Why are the following called the **success/failure** conditions?

p = proportion of "successes"

$1-p$ = proportion of "failures"

$$np \geq 10$$

$$n(1-p) \geq 10$$

Distribution of \hat{p}

Why are the following called the **success/failure** conditions?

p = proportion of "successes"
 $1-p$ = proportion of "failures"

$np \geq 10$
 $n(1-p) \geq 10$

Number of successes ←
 Number of failures ←

Distribution of \hat{p}

Central Limit Theorem for Proportions
 When **certain conditions** are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

Conditions:

- ▶ **Independence:**
 - a. Random sample/assignment AND
 - b. $n < 10\%$ of population
- ▶ **Shape/Skew:**
 - a. **"Success/Failure Conditions:"**
 - i. $np \geq 10$ (ie: at least 10 successes) AND
 - ii. $n(1-p) \geq 10$ (ie: at least 10 failures)

Clicker question

Suppose $p = 0.05$. What shape does the distribution of \hat{p} have in random samples of $n = 100$.

- (a) unimodal and symmetric (nearly normal)
- (b) bimodal and symmetric
- (c) right skewed
- (d) left skewed

Clicker question

Suppose $p = 0.05$. What shape does the distribution of \hat{p} have in random samples of $n = 100$.

~~(a) unimodal and symmetric (nearly normal)~~ **p and n don't satisfy SF Conditions.**

~~(b) bimodal and symmetric~~ **$np = 5 < 10$**

~~(c) right skewed~~ **$n(1-p) = 95 \geq 10$**

(d) left skewed

Clicker question



Suppose $p = 0.05$. What shape does the distribution of \hat{p} have in random samples of $n = 100$.

- (a) ~~unimodal and symmetric (nearly normal)~~ **p and n don't satisfy SF Conditions.** $np = 5 < 10$
- (b) ~~bimodal and symmetric~~ $n(1 - p) = 95 \geq 10$
- (c) *right skewed*
- (d) left skewed

If p and n don't satisfy SF conditions:

- Not unimodal and symmetric
- Sampling distribution for \hat{p} has one peak at p
- think about boundaries for \hat{p} , ie $[0,1]$.

Clicker question



Suppose $p = 0.5$. What shape does the distribution of \hat{p} have in random samples of $n = 100$.

- (a) unimodal and symmetric (nearly normal)
- (b) bimodal and symmetric
- (c) right skewed
- (d) left skewed

Clicker question



Suppose $p = 0.5$. What shape does the distribution of \hat{p} have in random samples of $n = 100$.

- (a) *unimodal and symmetric (nearly normal)* **p and n satisfy SF Conditions.** $np = 50 \geq 10$
- (b) bimodal and symmetric $n(1 - p) = 50 \geq 10$
- (c) right skewed
- (d) left skewed

Outline

1. Housekeeping

2. Main ideas

Hypothesis Testing and Confidence Intervals with a Single Categorical Variable

1. Theoretical Hypothesis Testing and Confidence Intervals: The CLT also describes the distribution of \hat{p}
2. Interpretations: CI vs. HT determines observed vs. expected counts / proportions
3. CLT Conditions: Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution

3. Applications

1. Single population proportion, large sample
2. Single population proportion, small sample

4. Recap

5. Summary

CI vs. HT determines observed vs. expected counts / proportions

Central Limit Theorem for Proportions
When certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

For instance...

$$P(\hat{p} < \text{obs}) = P(Z < \frac{\text{obs} - \text{mean}}{\text{stand.dev}}) = \text{probability}$$

Use z-tables

CI vs. HT determines observed vs. expected counts / proportions

Central Limit Theorem for Proportions
When certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

For instance...

$$P(\hat{p} < \text{obs}) = P(Z < \frac{\text{obs} - \text{mean}}{\text{stand.dev}}) = \text{probability}$$

Use z-tables

CI vs. HT determines observed vs. expected counts / proportions

Central Limit Theorem for Proportions
When certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

(1-α)% Confidence Interval
 $(\text{point estimate}) \pm z_{\alpha/2}^* SE$

Test Statistic for Hypothesis Test
 $z = \frac{(\text{point estimate}) - (\text{null value})}{SE}$

CI vs. HT determines observed vs. expected counts / proportions

Central Limit Theorem for Proportions
When certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

(1-α)% Confidence Interval
 $(\text{point estimate}) \pm z_{\alpha/2}^* SE$

Test Statistic for Hypothesis Test
 $z = \frac{(\text{point estimate}) - (\text{null value})}{SE}$

CI vs. HT determines observed vs. expected counts / proportions

Central Limit Theorem for Proportions
When certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev/error} = \sqrt{\frac{p(1-p)}{n}})$$

(1- α)% Confidence Interval
(point estimate) $\pm z_{\alpha/2}^* SE$

Test Statistic for Hypothesis Test
$$z = \frac{(\text{point estimate}) - (\text{null value})}{SE}$$

Problem: However, in Confidence intervals and Hypothesis Testing, we don't know what p is....

CI vs. HT determines observed vs. expected counts / proportions

Problem: In Confidence intervals and Hypothesis Testing, we don't know what p is but we need to plug it in in two parts of the analyses.

What part of the analyses do we need to plug in p ?	Confidence Interval for p	Hypothesis Test for p

CI vs. HT determines observed vs. expected counts / proportions

Problem: In Confidence intervals and Hypothesis Testing, we don't know what p is but we need to plug it in in two parts of the analyses.

What part of the analyses do we need to plug in p ?	Confidence Interval for p	Hypothesis Test for p
Checking SF Conditions: $np \geq 10$ $n(1-p) \geq 10$		
Calculating Standard Error: $SE = \sqrt{\frac{p(1-p)}{n}}$		

CI vs. HT determines observed vs. expected counts / proportions

Problem: In Confidence intervals and Hypothesis Testing, we don't know what p is but we need to plug it in in two parts of the analyses.

What part of the analyses do we need to plug in p ?	Confidence Interval for p	Hypothesis Test for p
Checking SF Conditions: $np \geq 10$ $n(1-p) \geq 10$	Solution: Plug in the observed \hat{p} for p .	
Calculating Standard Error: $SE = \sqrt{\frac{p(1-p)}{n}}$	Solution: Plug in the observed \hat{p} for p .	

CI vs. HT determines observed vs. expected counts / proportions

Problem: In Confidence intervals and Hypothesis Testing, we don't know what p is but we need to plug it in in two parts of the analyses.

What part of the analyses do we need to plug in p ?	Confidence Interval for p	Hypothesis Test for p
Checking SF Conditions: $np \geq 10$ $n(1-p) \geq 10$	Solution: Plug in the observed \hat{p} for p .	Solution: Plug in the null value p_0 for p .
Calculating Standard Error: $SE = \sqrt{\frac{p(1-p)}{n}}$	Solution: Plug in the observed \hat{p} for p .	Solution: Plug in the null value p_0 for p .

CI vs. HT determines observed vs. expected counts / proportions

Answer:

Remember, when doing a HT always assume H_0 is true!

- **S-F:** Number of successes and failures for checking the success-failure condition for the nearly normal distribution of \hat{p} :
 - CI: use observed proportion $\rightarrow n\hat{p} \geq 10$ and $n(1-\hat{p}) \geq 10$
 - HT: use null value of the proportion $\rightarrow np_0 \geq 10$ and $n(1-p_0) \geq 10$
- **SE:** Proportion of success for calculating the standard error of \hat{p} :

$$SE = \sqrt{\frac{p(1-p)}{n}}$$
 - CI: use observed proportion $\rightarrow SE = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$
 - HT: use null value of the proportion $\rightarrow SE = \sqrt{\frac{p_0(1-p_0)}{n}}$

NEW CI vs. HT determines observed vs. expected counts / proportions

Interpreting...

np = expected number of *successes*
 $n(1-p)$ = expected number of *failures*


$n\hat{p}$ = observed number of *successes*
 $n(1-\hat{p})$ = observed number *failures*

np_0 = expected number of *successes* (assuming H_0 is true)
 $n(1-p_0)$ = expected number of *failures* (assuming H_0 is true)

Outline

- Housekeeping
- Main ideas
 - Hypothesis Testing and Confidence Intervals with a Single Categorical Variable**
 - Theoretical Hypothesis Testing and Confidence Intervals: The CLT also describes the distribution of \hat{p}
 - Interpretations: CI vs. HT determines observed vs. expected counts / proportions
 - CLT Conditions: Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution
- Applications
 - Single population proportion, large sample
 - Single population proportion, small sample
- Recap
- Summary

Single Population Proportion, p



When to use each method

- Construct confidence intervals and hypothesis tests using **Central Limit Theorem methods**.
- Conduct a hypothesis test for p with **randomization testing**.
- Create a confidence interval for p with **bootstrapping**.

Simulation vs. theoretical inference

Confidence Interval

Conditions:

- ▶ **Independence:**
 - a. Random sample/assignment
 - b. 10% rule
- ▶ **Shape/Skew:**
 - a. "Success/Failure Conditions:" At least 10 successes and failures (ie: $n\hat{p} \geq 10, n(1-\hat{p}) \geq 10$)

Condition met: Condition NOT met:

$(1-\alpha)\%$ Confidence Interval
With CLT Methods

$$\hat{p} \pm z_{\alpha/2}^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Make a **bootstrap confidence interval.**

Simulation vs. theoretical inference

Hypothesis Test

Conditions:

- ▶ **Independence:**
 - a. Random sample/assignment
 - b. 10% rule
- ▶ **Shape/Skew:**
 - a. "Success/Failure Conditions:" At least 10 successes and failures (ie: $np_0 \geq 10, n(1-p_0) \geq 10$)

Condition met: Condition NOT met:

CLT Methods
Test Statistic

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

Conduct hypothesis test for p with **randomization testing** – balls/chips in a bag, dice, etc.

Simulation vs. theoretical inference

Central Limit Theorem for Proportions

When certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev/error} = \sqrt{\frac{p(1-p)}{n}})$$

$(1-\alpha)\%$ Margin of Error
 $ME = z_{\alpha/2}^* SE$

Problem: In Margin of Error calculation:

- We don't know what p is.
- We may also not know what \hat{p} is.

Recap on CLT based methods

► **Calculating the necessary sample size for a CI with a given margin of error:**

– Option 1: If there is a previous study, use \hat{p} from that study

Recap on CLT based methods

► **Calculating the necessary sample size for a CI with a given margin of error:**

– Option 1: If there is a previous study, use \hat{p} from that study

– Option 2: If not, use $\hat{p} = 0.5$:

- if you don't know any better, 50-50 is a good guess
- $\hat{p} = 0.5$ gives the most conservative estimate of the standard error, which gives the highest possible sample size

Outline

1. Housekeeping

2. Main ideas

Hypothesis Testing and Confidence Intervals with a Single **Categorical** Variable

1. Theoretical Hypothesis Testing and Confidence Intervals: \mathcal{Q} The CLT also describes the distribution of \hat{p}
2. Interpretations: \mathcal{H} CI vs. HT determines observed vs. expected counts / proportions
3. CLT Conditions: \mathcal{Q} Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution

3. Applications

1. \mathcal{H} Single population proportion, large sample
2. \mathcal{H} \mathcal{H} Single population proportion, small sample

4. Recap

5. Summary

Application exercise: App Ex 5.1

See course website for details.



Outline

1. Housekeeping

2. Main ideas

Hypothesis Testing and Confidence Intervals with a Single **Categorical** Variable

1. [Theoretical Hypothesis Testing and Confidence Intervals](#): \mathcal{Q} The CLT also describes the distribution of \hat{p}
2. [Interpretations](#): \mathcal{NEW} \mathcal{CI} vs. HT determines observed vs. expected counts / proportions
3. [CLT Conditions](#): \mathcal{Q} Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution

3. Applications

1. \mathcal{H} Single population proportion, large sample
2. \mathcal{H} \mathcal{NEW} \mathcal{CI} Single population proportion, small sample

4. Recap

5. Summary

Outline

\mathcal{H} Let's conduct a hypothesis test on p (the proportion of Duke students that are vegetarian/vegan) with **randomization testing**.



Clicker question



Are you vegetarian or vegan?

- (a) Yes, I am vegetarian or vegan
- (b) No, I am neither vegetarian nor vegan





Clicker question



A variety of studies suggest that 8% of college students are vegetarian or vegan. Assuming that this class is a representative sample of Duke students, which of the following are the correct set of hypotheses for testing if the proportion of Duke students who are vegetarian is different than the proportion of vegetarian college students at large.



- (a) $H_0 : p = 0.08; H_A : p \neq 0.08$
- (b) $H_0 : p = 0.08; H_A : p < 0.08$
- (c) $H_0 : \hat{p} = 0.08; H_A : \hat{p} \neq 0.08$
- (d) $H_0 : \hat{p}_{Duke} = \hat{p}_{all\ college}; H_A : \hat{p}_{Duke} \neq \hat{p}_{all\ college}$
- (e) $H_0 : p_{Duke} = p_{all\ college}; H_A : p_{Duke} \neq p_{all\ college}$

Clicker question  

A variety of studies suggest that 8% of college students are vegetarian or vegan. Assuming that this class is a representative sample of Duke students, which of the following are the correct set of hypotheses for testing if the proportion of Duke students who are vegetarian is **different** than the proportion of vegetarian college students at large.



(a) $H_0 : p = 0.08; H_A : p \neq 0.08$
 (b) $H_0 : p = 0.08; H_A : p < 0.08$
 (c) $H_0 : \hat{p} = 0.08; H_A : \hat{p} \neq 0.08$
 (d) $H_0 : \hat{p}_{Duke} = \hat{p}_{all\ college}; H_A : \hat{p}_{Duke} \neq \hat{p}_{all\ college}$
 (e) $H_0 : p_{Duke} = p_{all\ college}; H_A : p_{Duke} \neq p_{all\ college}$

"A variety of studies..." usually suggests we can assume that we "know" the pop. parameter it's referring to (0.08).

  Simulate by hand

Goal in Randomization Testing:
 1. Set up hypotheses.



Example 1:
 $H_0: p=0.08 (=8/100)$
 $H_a: p \neq 0.08$

  Simulate by hand

Goal in Randomization Testing:
 1. Set up hypotheses.
 2. Create a Randomization Distribution: an approximation for the sampling distribution that **assumes H_0 is true (or is centered at the null value)**

Example 1:
 $H_0: p=0.08 (=8/100)$
 $H_a: p \neq 0.08$

- ▶ 100 chips in a bag: 8 green (veg), 92 white (non veg).
- ▶ Sample randomly n times from the bag, **with replacement** ($n =$ observed sample size=144)
- ▶ Calculate \hat{p} , the proportion of greens (successes) in the random sample of size $n=144$, record this value.
- ▶ Repeat many times.

  Simulate by hand

Goal in Randomization Testing:
 1. Set up hypotheses.
 2. Create a Randomization Distribution: an approximation for the sampling distribution that **assumes H_0 is true (or is centered at the null value)**
 3. Calculate p-value with this randomization distribution: **P-value**=% of simulated points in the randomization distribution that are **at least as extreme** as the sample statistic observed.

Example 1:
 $H_0: p=0.08 (=8/100)$
 $H_a: p \neq 0.08$

- ▶ 100 chips in a bag: 8 green (veg), 92 white (non veg).
- ▶ Sample randomly n times from the bag, **with replacement** ($n =$ observed sample size=144)
- ▶ Calculate \hat{p} , the proportion of greens (successes) in the random sample of size $n=144$, record this value.
- ▶ Repeat many times.
- ▶ Calculate the proportion of simulations where \hat{p} is **at least as different from 0.08 (null value) as the observed sample proportion.**

Clicker question NEW

What would we expect this randomization distribution to be centered at?

- a) 0
- b) 0.08
- c) the proportion of vegans/vegetarians in this class
- d) Some number higher than 0.08.

Clicker question NEW

What would we expect this randomization distribution to be centered at?

- a) 0
- b) **0.08 – The null value.**
- c) the proportion of vegans/vegetarians in this class
- d) Some number higher than 0.08.

NEW Simulate by hand

Ho: $p=0.08$
 Ha: $p \neq 0.08$

Example:

P-value=

50 simulated sample proportions

Simulation Proportions

NEW HT in R

```

n_veg = # of veg's in class
n_nonveg = # of nonveg's in class
sta101 = data.frame(veg = c(rep("yes", n_veg), rep("no", n_nonveg)))

inference(y = veg, data = sta101, success = "yes",
          statistic = "proportion", type = "ht",
          null = 0.08, alternative = "twosided",
          method = "simulation")
    
```

Outline

👋 Let's create a **bootstrap confidence interval** for p (the proportion of Duke students that are vegetarian/vegan).



Bootstrap interval for a single proportion

How would the simulation scheme change for a bootstrap interval for the proportion of Duke students who are vegetarians?



Bootstrap interval for a single proportion

How would the simulation scheme change for a bootstrap confidence interval for the proportion of Duke students who are vegetarians?

Step 1: Create a Bootstrap Distribution

- ▶ ~~100 chips in a bag: 8 green (veg), 92 white (non veg).~~
- ▶ Sample randomly n times from the ~~bag~~ **original sample**, with replacement ($n =$ observed sample size)
- ▶ Calculate \hat{p} , the proportion of **vegetarians** (successes) in the random sample of size n , record this value.
- ▶ Repeat many times.



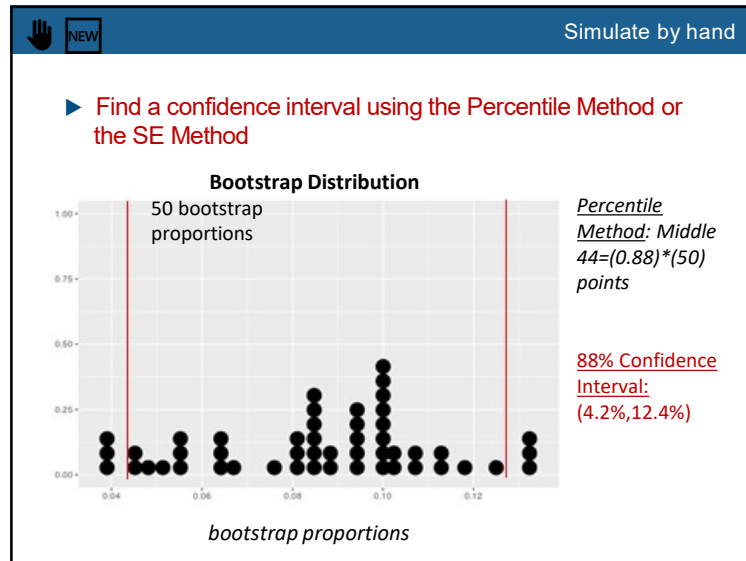
Bootstrap interval for a single proportion

How would the simulation scheme change for a bootstrap confidence interval for the proportion of Duke students who are vegetarians?

Step 1: Create a Bootstrap Distribution

- ▶ ~~100 chips in a bag: 8 green (veg), 92 white (non veg).~~
- ▶ Sample randomly n times from the ~~bag~~ **original sample**, with replacement ($n =$ observed sample size)
- ▶ Calculate \hat{p} , the proportion of **vegetarians** (successes) in the random sample of size n , record this value.
- ▶ Repeat many times.

Step 2: Using the Bootstrap distribution, find a confidence interval using the Percentile Method or the SE Method



CI in R

```
inference(y = veg, data = sta101, success = "yes",
          statistic = "proportion", type = "ci",
          method = "simulation", boot_method = "se")
```

- Summary of main ideas
1. The CLT also describes the distribution of \hat{p}
 2. CI vs. HT determines observed vs. expected counts / proportions
 3. Only use CLT based methods if the sample size is large enough for a nearly normal sampling distribution