# Hypothesis Tests and Confidence Intervals using Normal Distribution!

STAT 120

Sections 5.1-5.2

Day 15

# A Research Question

## How do Malaria parasites impact mosquito behavior?

# Malaria Parasites and Mosquitoes

- Mosquitos were randomized to either eat from a malaria infected mouse (exposed group) or a healthy mouse (control group)

- After infection, the parasites go through two stages:

  1) Not yet infectious (Days 1-8)

  2) Infectious (Days 9 – 28)

- Response variable: whether the mosquito approached a human (in a cage with them)

- Does this behavior differ by exposed vs control?  Does it differ by infection stage?

Cator LJ, George J, Blanford S, Murdock CC, Baker TC, Read AF, Thomas MB. (2013). 'Manipulation' without the parasite: altered feeding behaviour of mosquitoes is not dependent on infection with malaria parasites. Proc R Soc B 280: 20130711.

## Malaria Parasites and Mosquitoes

- Does infecting mosquitoes with Malaria actually impact the mosquitoes' behavior to favor the parasite?

- Malaria parasites would benefit if:

  - Mosquitoes approached humans less often after being exposed, but before becoming infectious, because humans are risky

  - Mosquitoes approached humans more often after becoming infectious, to pass on the infection

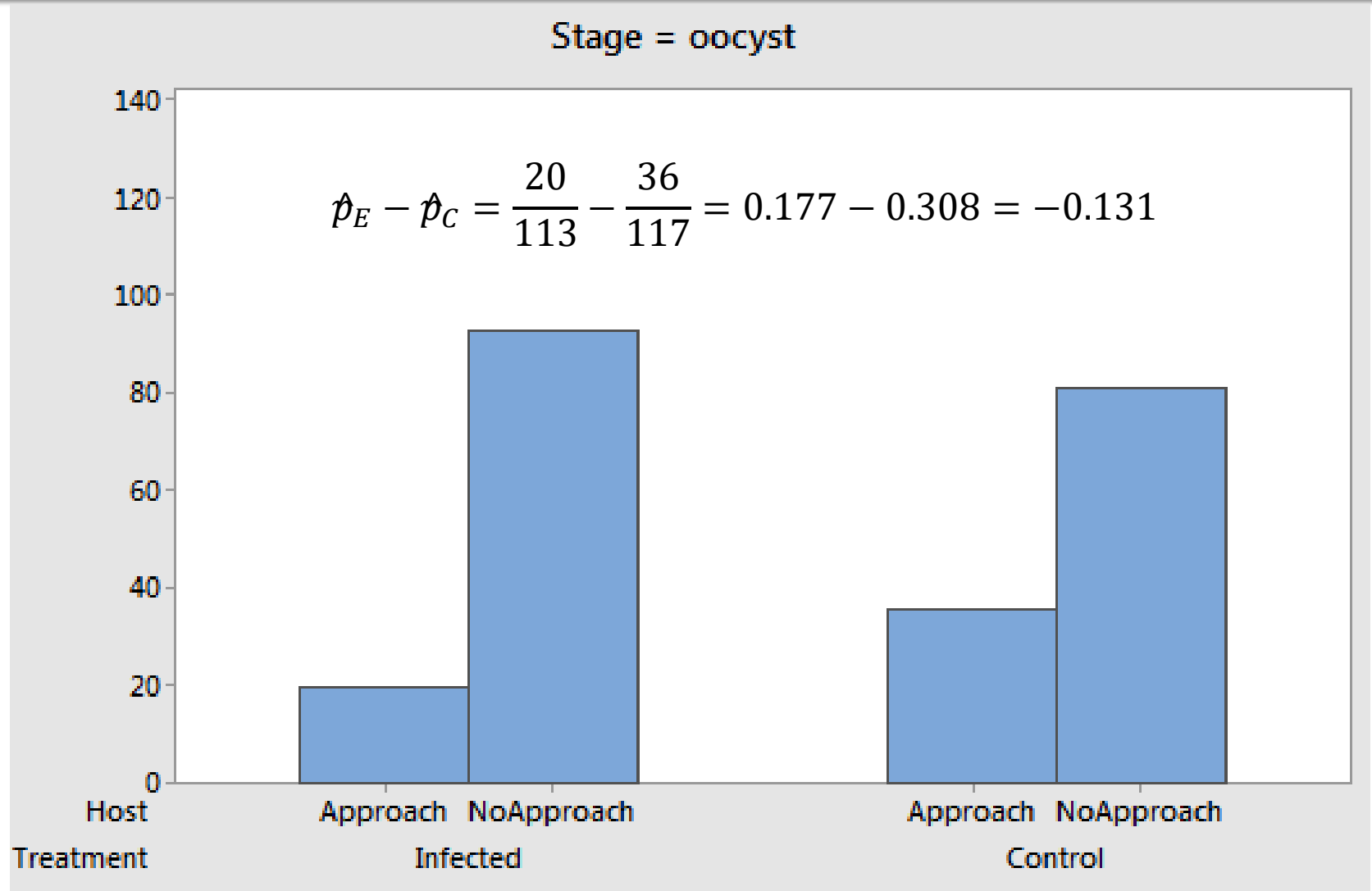**We'll first look at the mosquitoes before they become infectious (days 1–8).**

$p_C$: proportion of controls to approach human

$p_E$: proportion of exposed to approach human
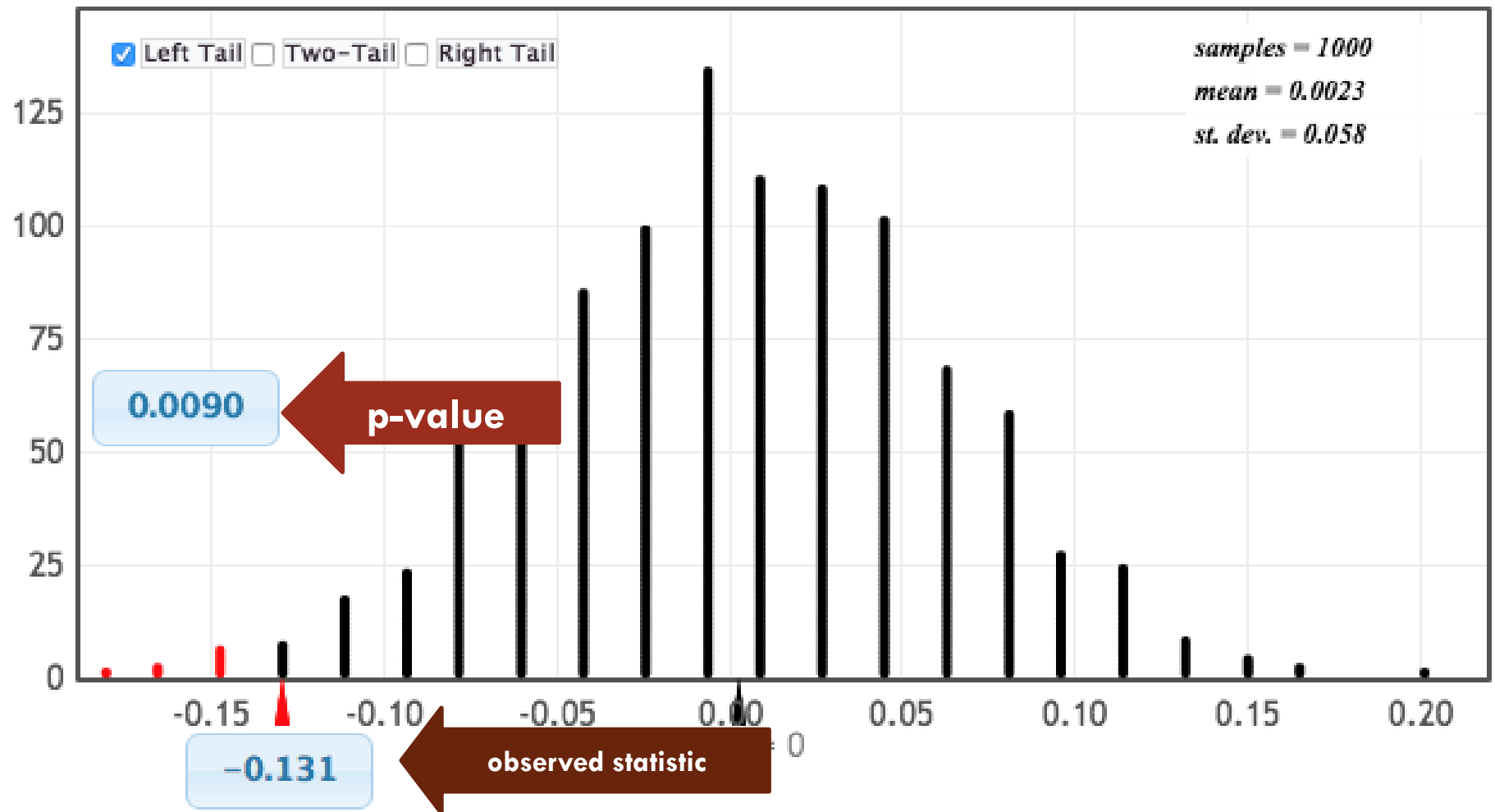
**What are the relevant hypotheses?**

A. $H_0: p_E = p_C$, $H_a: p_E < p_C$

B. $H_0: p_E = p_C$, $H_a: p_E > p_C$

C. $H_0: p_E < p_C$, $H_a: p_E = p_C$

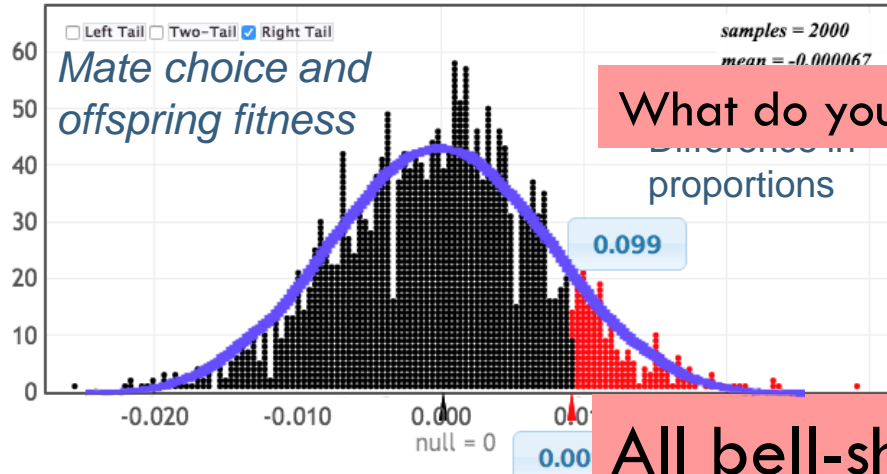D. $H_0: p_E > p_C$, $H_a: p_E = p_C$

# Data: Before Infectious



Stage = oocyst

$$\hat{p}_E - \hat{p}_C = \frac{20}{113} - \frac{36}{117} = 0.177 - 0.308 = -0.131$$

# Randomization Test



Randomization Dotplot of $\hat{p}_1 - \hat{p}_2$ ▼    Null Hypothesis: $p_1 = p_2$

Left Tail ☐ Two-Tail ☐ Right Tail

samples = 1000
mean = 0.0023
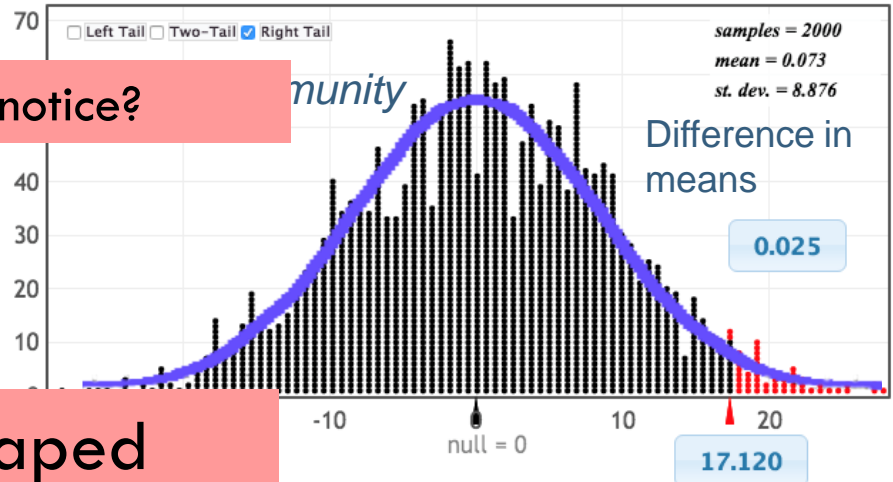st. dev. = 0.058

0.0090  p-value

−0.131  observed statistic

# Randomization and Bootstrap Distributions



**Randomization Dotplot of** $\hat{p}_1 - \hat{p}_2$   **Null Hypothesis:** $p_1 = p_2$

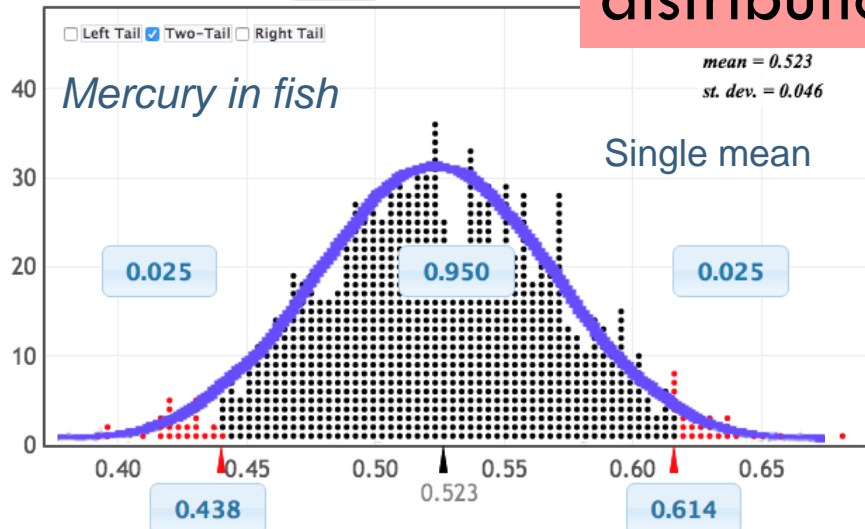*Mate choice and offspring fitness*

Left Tail ☐  Two-Tail ☐  Right Tail ☑

samples = 2000
mean = -0.000067

Difference in proportions

0.099

null = 0

0.00

**Randomization Dotplot of** $\bar{x}_1 - \bar{x}_2$ ,  **Null hypothesis:** $\mu_1 = \mu_2$

...munity

Left Tail ☐  Two-Tail ☐  Right Tail ☑

samples = 2000
mean = 0.073
st. dev. = 8.876

Difference in means

0.025

null = 0

17.120

**Bootstrap Dotplot of** Mean ▾

*Mercury in fish*

Left Tail ☐  Two-Tail ☑  Right Tail ☐

mean = 0.523
st. dev. = 0.046

Single mean

0.025   0.950   0.025

0.438   0.523   0.614

...lot of $\bar{x}_1 - \bar{x}_2$ ,  **Null hypothesis:** $\mu_1 = \mu_2$

*Sleep versus caffeine*

Left Tail ☐  Two-Tail ☑  Right Tail ☐

samples = 1000
mean = -0.011
st. dev. = 1.491

Difference in proportions

0.024   0.952   0.024

null = 0

−3.000   3.000

What do you notice?

All bell-shaped distributions!

# Central Limit Theorem

For random samples with a sufficiently large sample size, the distribution of sample statistics for a **mean** or a **proportion** is **normally distributed**
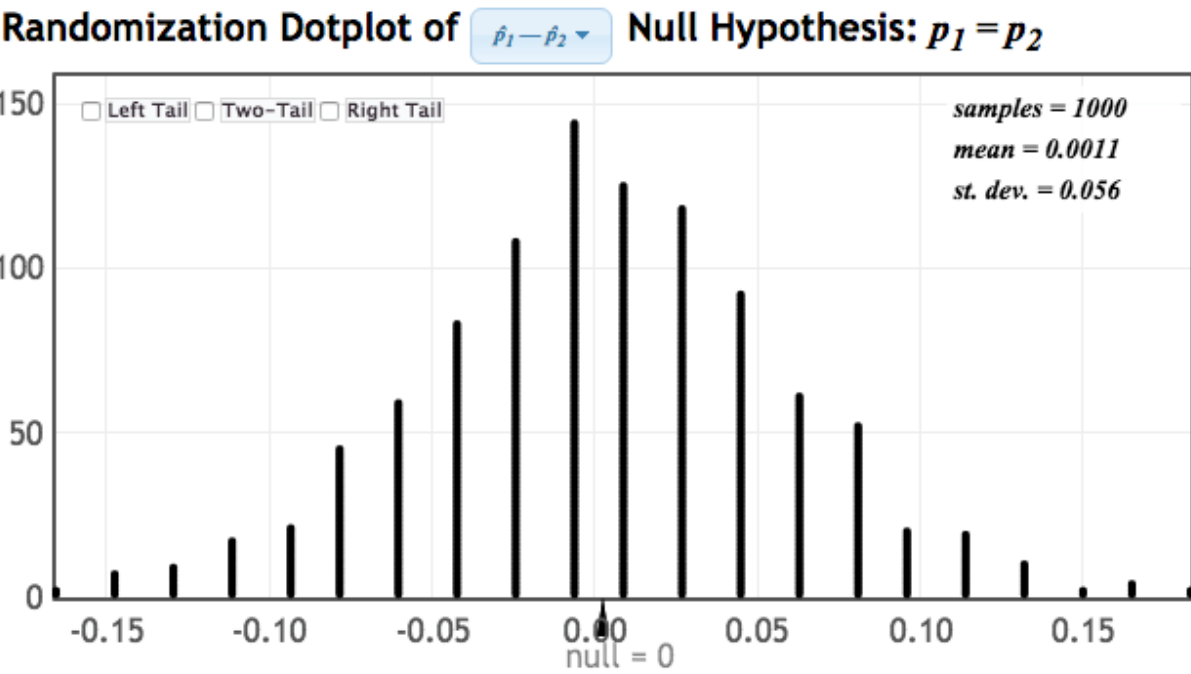
# Central Limit Theorem

- The catch:

"sufficiently large sample size"

- The **more skewed** the original distribution of data/population is, the larger $n$ has to be for the CLT to work

- For quantitative variables that are not very skewed, $n \geq 30$ is usually sufficient

- For categorical variables, counts of at least 10 within each category is usually sufficient

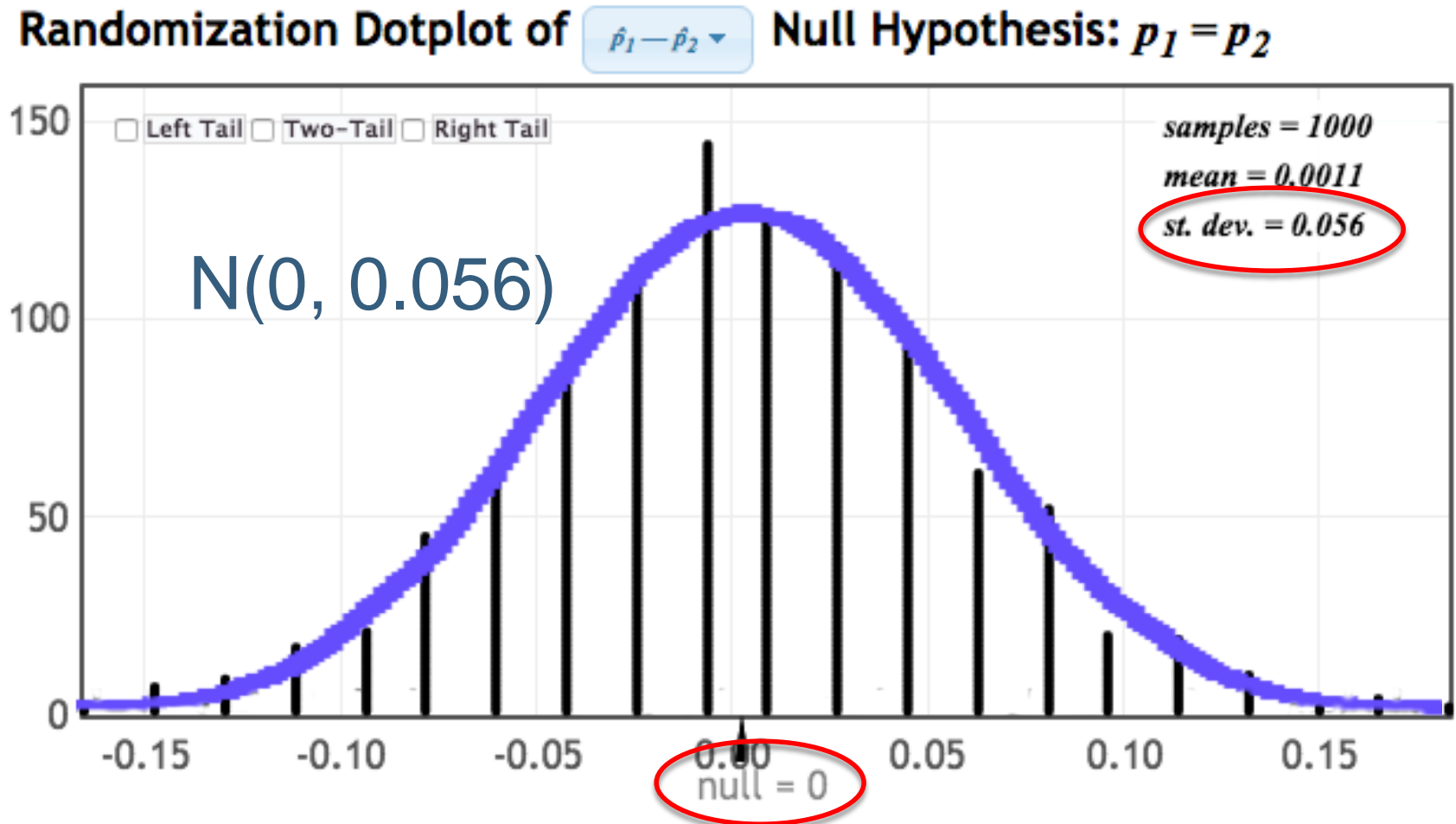# Which normal distribution should we use to approximate this?



Randomization Dotplot of $\hat{p}_1 - \hat{p}_2$ ▾   Null Hypothesis: $p_1 = p_2$

☐ Left Tail ☐ Two-Tail ☐ Right Tail

samples = 1000
mean = 0.0011
st. dev. = 0.056

null = 0

**Original Sample**

| Group | Count | Sample Size | Proportion |
|---|---|---|---|
| Group 1 | 20 | 113 | 0.177 |
| Group 2 | 36 | 117 | 0.308 |
| Group 1– Group 2 | –16 | n/a | –0.131 |

**Randomization Sample**

| Group | Count | Sample Size | Proportion |
|---|---|---|---|
| Group 1 | 27 | 113 | 0.239 |
| Group 2 | 29 | 117 | 0.248 |
| Group 1– Group 2 | –2 | n/a | –0.0089 |

A.  N(0, -0.131)

B.  N(0, 0.056)

C.  N(-0.131, 0.056)

D.  N(0.056, 0)

# Normal Distribution

**Randomization Dotplot of** $\hat{p}_1 - \hat{p}_2$ ▾ **Null Hypothesis:** $p_1 = p_2$



N(0, 0.056)

samples = 1000
mean = 0.0011
st. dev. = 0.056

☐ Left Tail ☐ Two–Tail ☐ Right Tail

null = 0

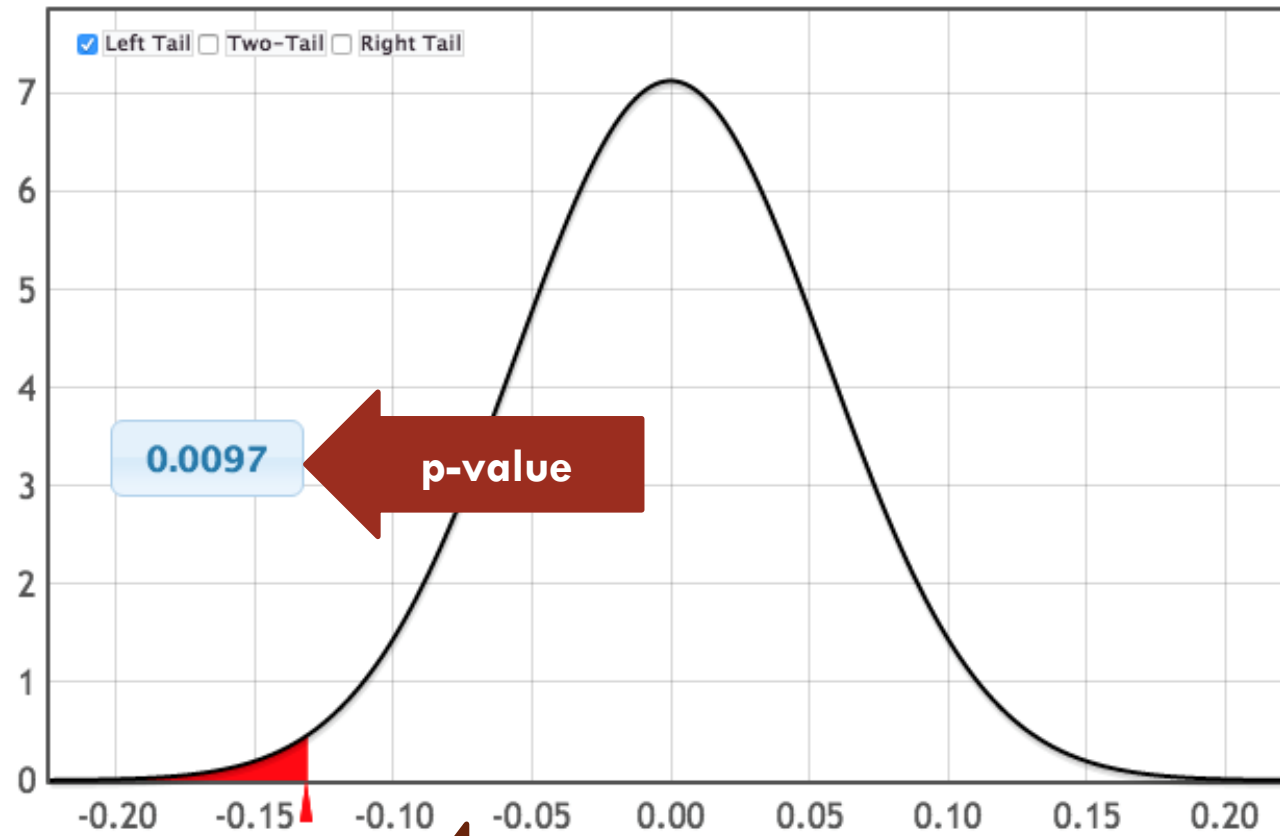We can compare the original statistic to this Normal distribution to find the p-value!

# p-value from N(null, SE)

# Connecting Normal model to hypothesis tests

- Hypothesis test:
  - Suppose: randomization distribution is bell-shaped.
  - Center: hypothesized null parameter value
  - Spread: the standard error given in the randomization graph (or by formula)
- P-value is computed from the normal model the "usual" way – the chance of being as extreme, or more extreme, than the observed statistic.
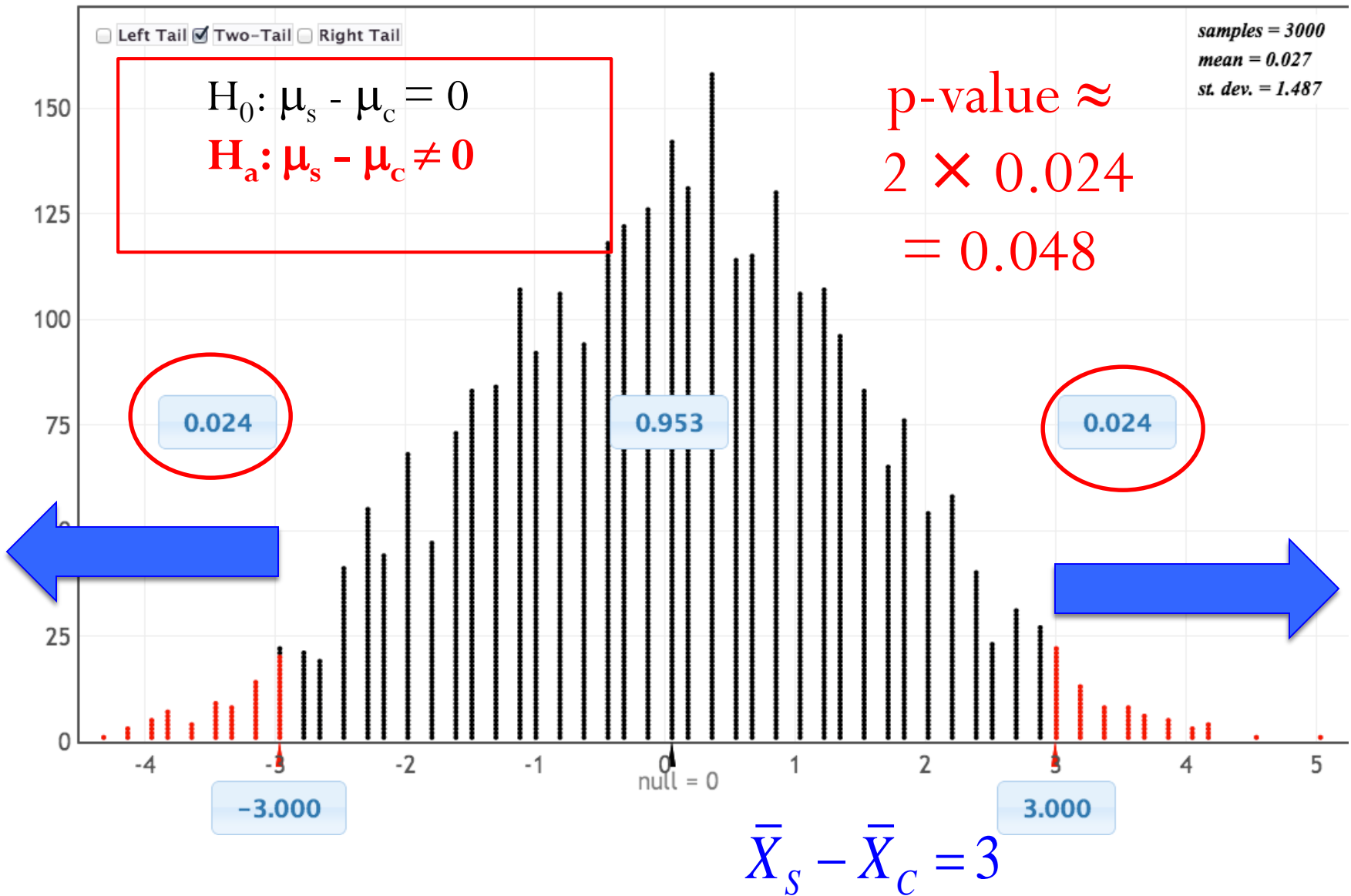
# Standardized Statistic

The **standardized test statistic** (also known as a z-statistic) is

$$z = \frac{statistic - null}{SE}$$

☐ Calculating the number of standard errors a statistic is from the null lets us assess extremity on a common scale.

# Sleep versus Caffeine

Randomization Dotplot of $\bar{x}_1 - \bar{x}_2$,  Null hypothesis: $\mu_1 = \mu_2$



Left Tail  ☑ Two-Tail  Right Tail

samples = 3000
mean = 0.027
st. dev. = 1.487

$H_0: \mu_s - \mu_c = 0$
$H_a: \mu_s - \mu_c \neq 0$

p-value $\approx$
$2 \times 0.024$
$= 0.048$

0.024     0.953     0.024

−3.000     3.000

null = 0

$\bar{X}_S - \bar{X}_C = 3$

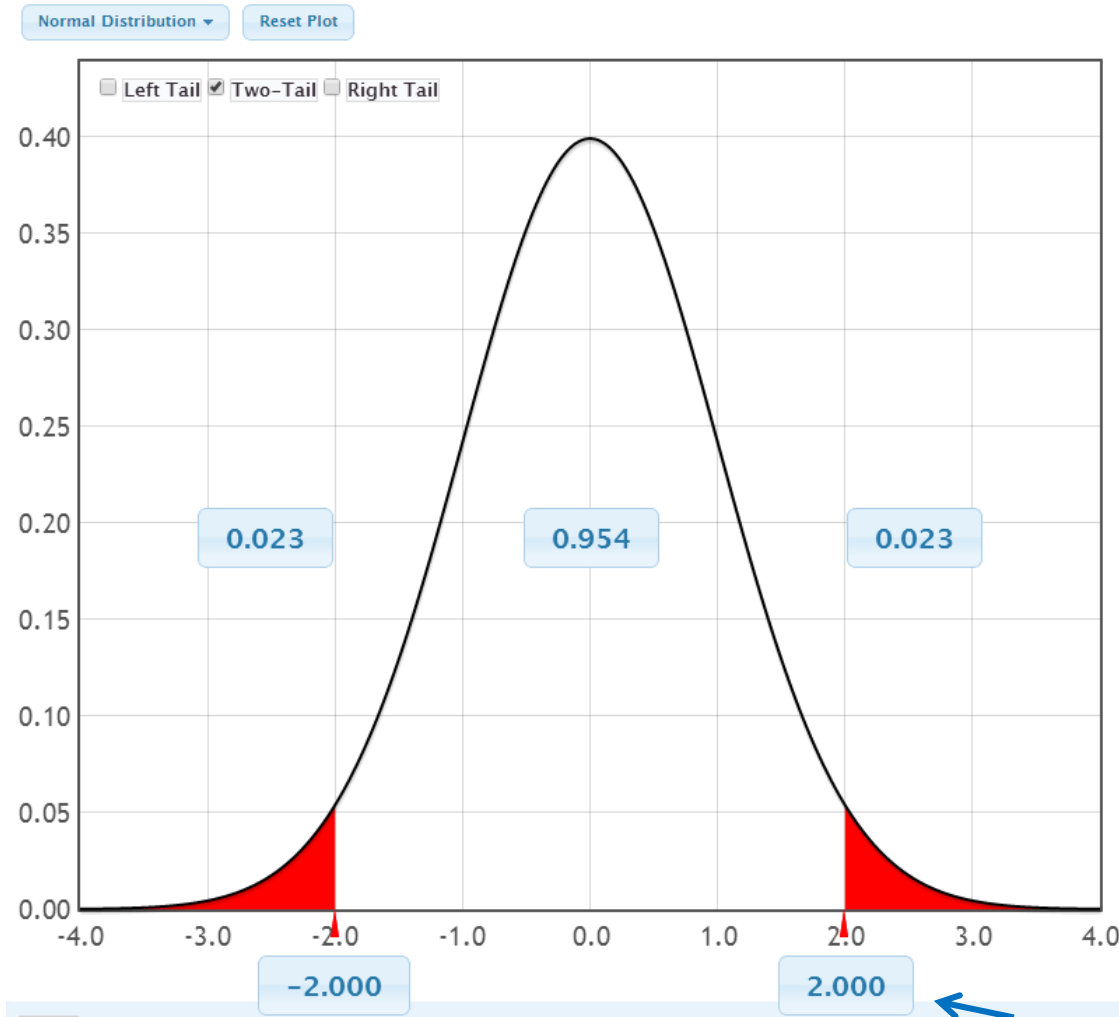# Sleep versus Caffeine

$$H_0: \mu_s - \mu_c = 0$$
$$H_a: \mu_s - \mu_c \neq 0$$

- The observed difference is 3 words.

- The randomization model is approximately $N(0, 1.5)$.

- The observed difference is 2 SEs above the null difference of 0:

$$z = \frac{\text{sample stat} - \text{null parameter value}}{\text{SE}} = \frac{3 - 0}{1.5} = 2$$

- The p-value is the area above +3 and below -3 on the $N(0, 1.5)$ model.

- Equivalently, the p-value is the chance of getting a mean difference more than 2 SEs away from the hypothesized mean difference of 0.

# Sleep versus Caffeine

**Normal Distribution ▾**  **Reset Plot**

☐ Left Tail ☑ Two-Tail ☐ Right Tail

**Normal Distribution**

| Mean | Standard Deviation |
|------|--------------------|
| 0 | 1 |

**Edit Parameters**

0.023   0.954   0.023

-2.000   2.000

**N(0,1) model!**

**The p-value is about 2(.023), or about 4.6%.**

```
> 2*pnorm(-2)
[1] 0.04550026
```

**Standardized Test Stat value**

# Malaria and Mosquitos

☐ Does infecting mosquitoes with Malaria actually impact the mosquitoes' behavior to favor the parasite?

☐ **After the parasite becomes infectious**, do infected mosquitoes approach humans *more* often, so as to pass on the infection?

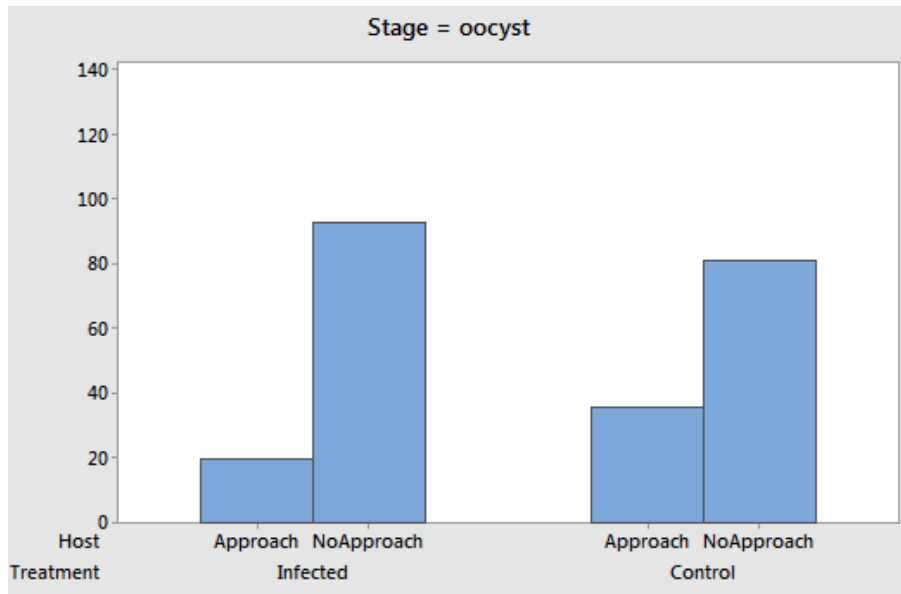**For the data after the mosquitoes become infectious (Days 9 – 28), what are the relevant hypotheses?**

$p_C$: proportion of controls to approach human
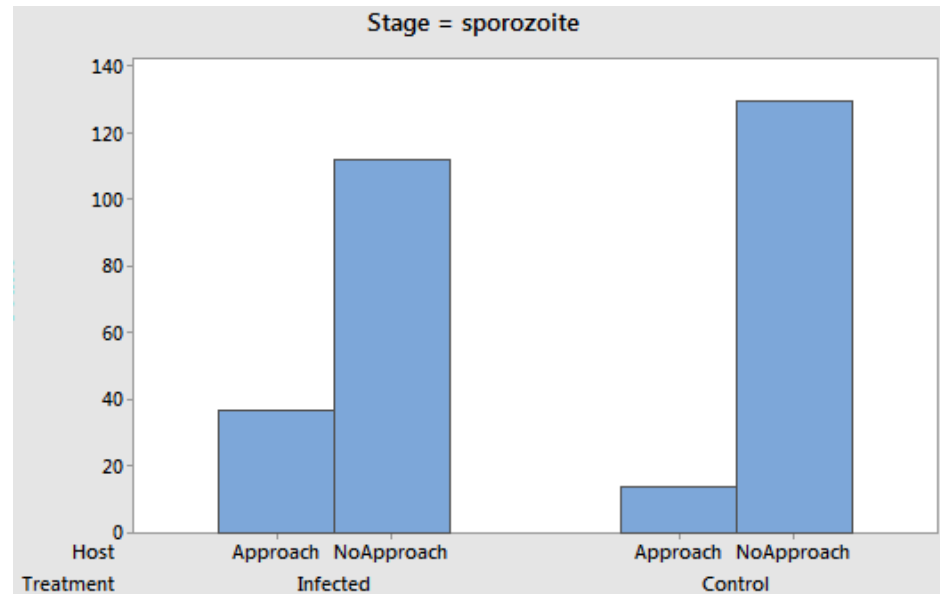
$p_E$: proportion of exposed to approach human

A. $H_0: p_E = p_C, H_a: p_E < p_C$

B. $H_0: p_E = p_C, H_a: p_E > p_C$

C. $H_0: p_E < p_C, H_a: p_E = p_C$

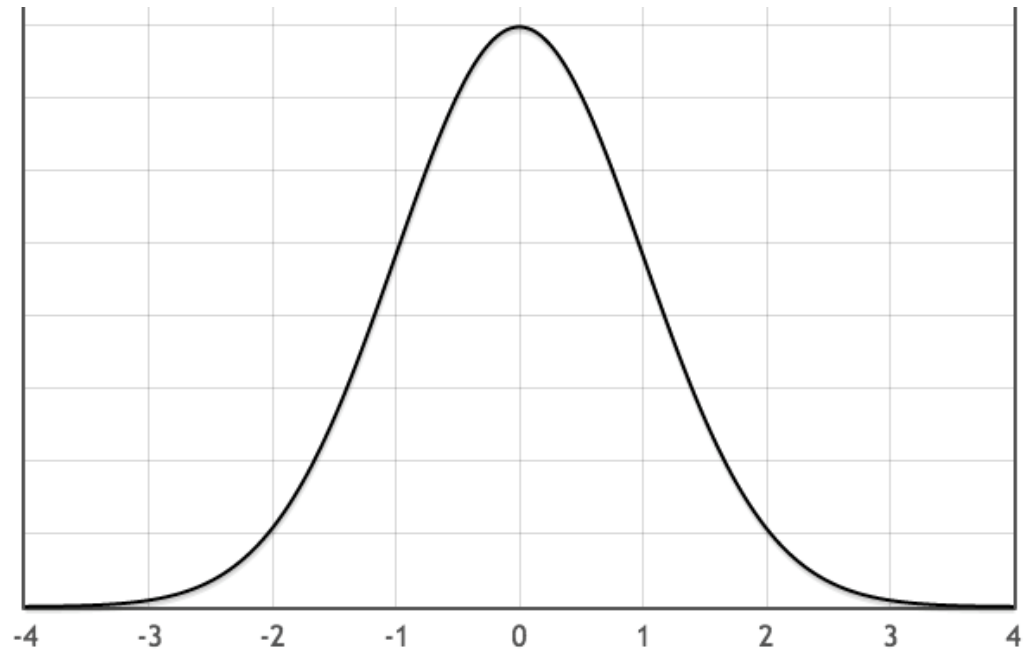D. $H_0: p_E > p_C, H_a: p_E = p_C$

# Data

## Before Infectious



## After Infectious

**The difference in proportions is 0.151 and the standard error is 0.05. Is this significant?**

A. Yes

B. No

# Malaria and Mosquitoes

- It appears that mosquitoes infected by malaria parasites do, in fact, behave in ways advantageous to the parasites!

  - Exposed mosquitos are *less* likely to approach before becoming infectious (so more likely to stay alive)

  - Exposed mosquitos are *more* likely to approach humans after becoming infectious (so more likely to pass on disease)

- Toxoplasmosis parasites impact rat and human behavior, malaria parasites impact mosquito behavior… what else might parasites be impacting that we have yet to discover??

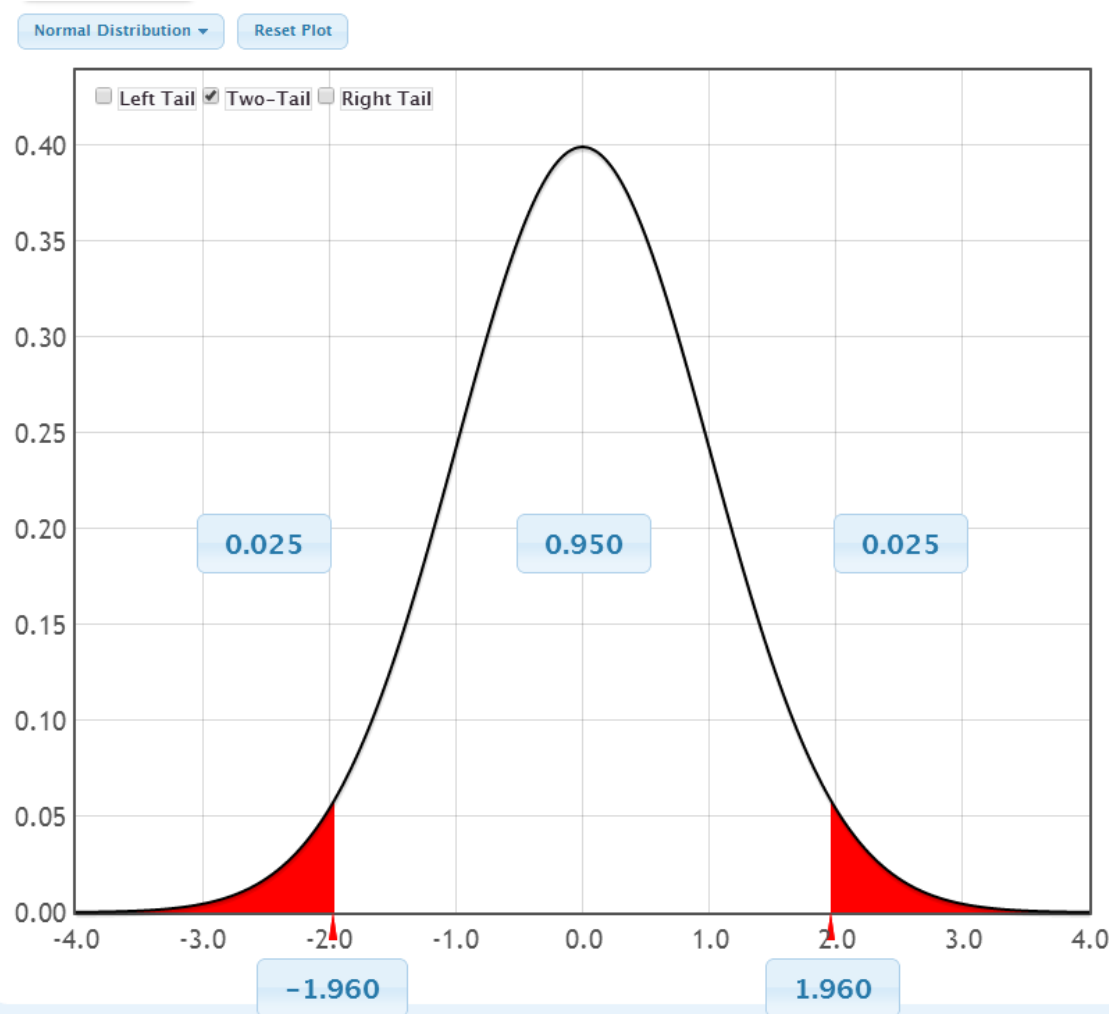# Formula for p-values Using N(0,1)

From original data

From $H_0$

$$z = \frac{\text{sample statistic} - \text{null value}}{\text{SE}}$$

From randomization distribution

# Connecting Normal model to Confidence Intervals

- Confidence Intervals
  - Suppose: bootstrap distribution is bell-shaped.
  - Center: sample statistic
  - Spread: the standard error given in the bootstrap graph (or by formula)
- To get a 95% confidence interval we compute

$$\textbf{statistic +/− 2(SE)}$$

- Why 2 SE's?
  - 95% of all sample means fall within 2 SE's of the population mean*
  - The value 2 is a z-score!
  - * well, actually the precise z-score under a normal model is z=1.96 instead of 2!
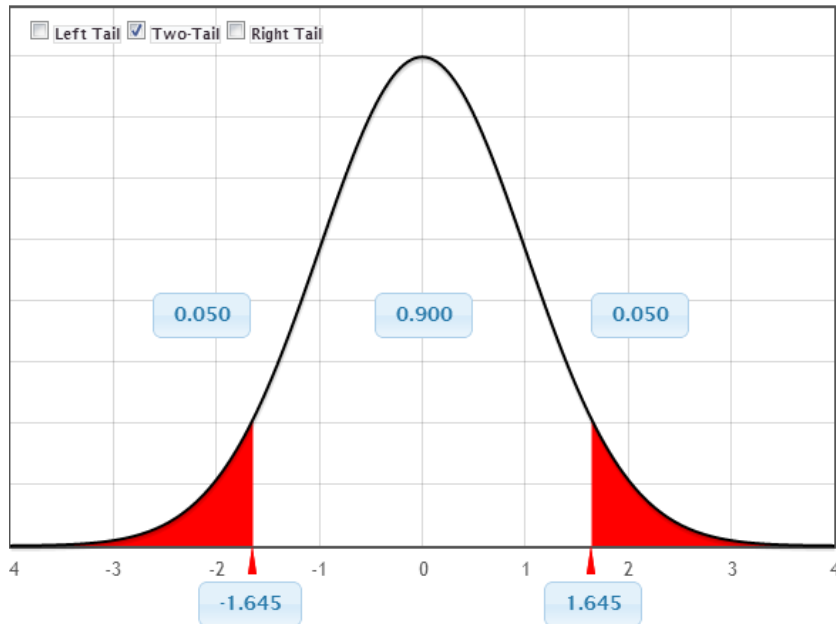
# N(0,1) model



95% of all values fall within **1.96** SE's of the mean

```
> qnorm(.975)
[1] 1.959964
```

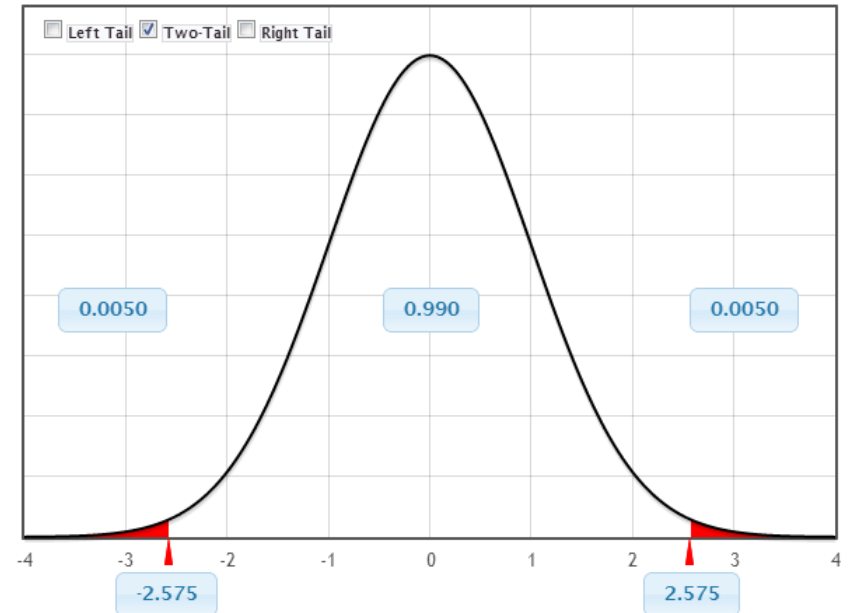What if we wanted a 90% CI? What z-score should we use to get the margin of error?

# Other Levels of Confidence

90% Confidence
$z* = 1.645$

```
> qnorm(.95)
[1] 1.644854
```

99% Confidence
$z* = 2.576$

```
> qnorm(.995)
[1] 2.575829
```

# Confidence Interval using N(0,1)

If a statistic is normally distributed, we find a confidence interval for the parameter using

*statistic +/− z\* SE*

where the area between −z\* and +z\* in the standard normal distribution is the desired level of confidence.

# Global Warming

What percentage of Americans believe in global warming?

A survey on 2,251 randomly selected individuals conducted in October 2010 found that 1328 answered "Yes" to the question

*"Is there solid evidence of global warming?"*

Give and interpret a 95% CI for the proportion of Americans who believe there is solid evidence of global warming.

# Global Warming

## Bootstrap For One Categorical Variable [Return to StatKey Index]
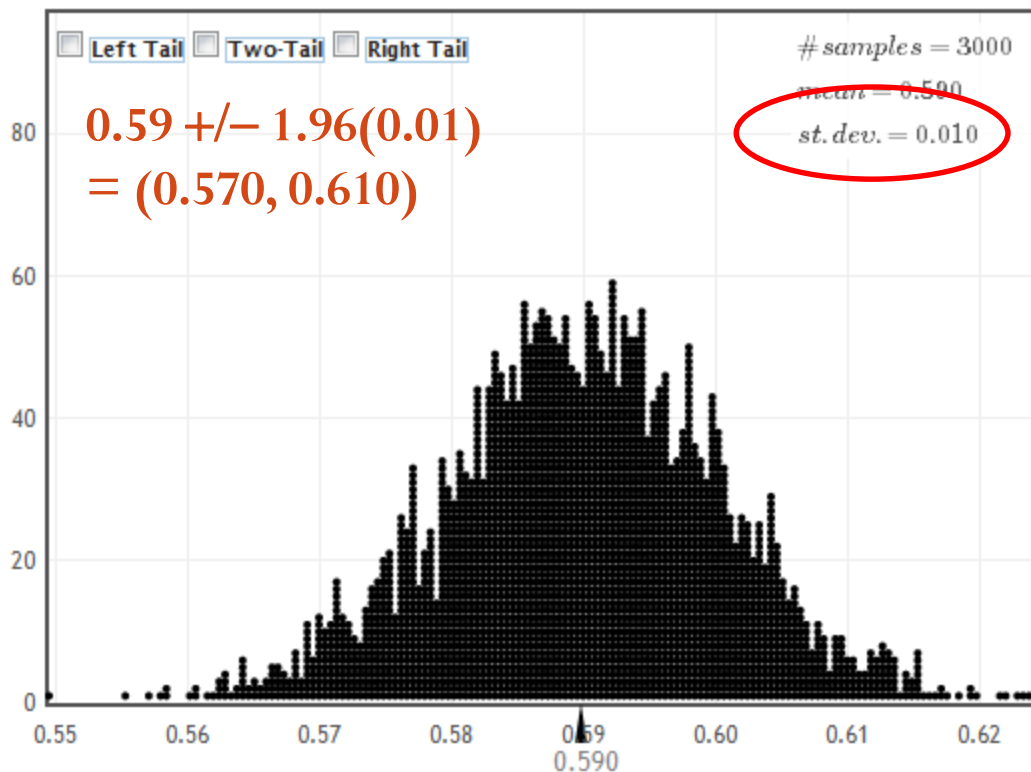
Custom Data ▼    Edit Data

Generate 1 Samples    Generate 10 Samples    Generate 100 Samples    Generate 1000 Samples    Reset Plot

**Bootstrap Dotplot of** Proportion ▼

☐ Left Tail  ☐ Two-Tail  ☐ Right Tail

$\#samples = 3000$

$mean = 0.590$

$st.dev. = 0.010$

**0.59 +/− 1.96(0.01)**
**= (0.570, 0.610)**

### Original Sample

| Count | n | Proportion |
|-------|------|------------|
| 1328 | 2251 | 0.590 |

### Bootstrap Sample

| Count | n | Proportion |
|-------|------|------------|
| 1304 | 2251 | 0.579 |

*We are 95% confident that the true percentage of all Americans that believe there is solid evidence of global warming is between 57.0% and 61.0%*

# Global Warming

- What is a 90% confidence interval for the proportion of US adults who believe in global warming?

**$z^*=1.645$ for 90% confidence**

**$0.59 +/- 1.645(0.01) = (0.574, 0.606)$**

- What is a 99% confidence interval?
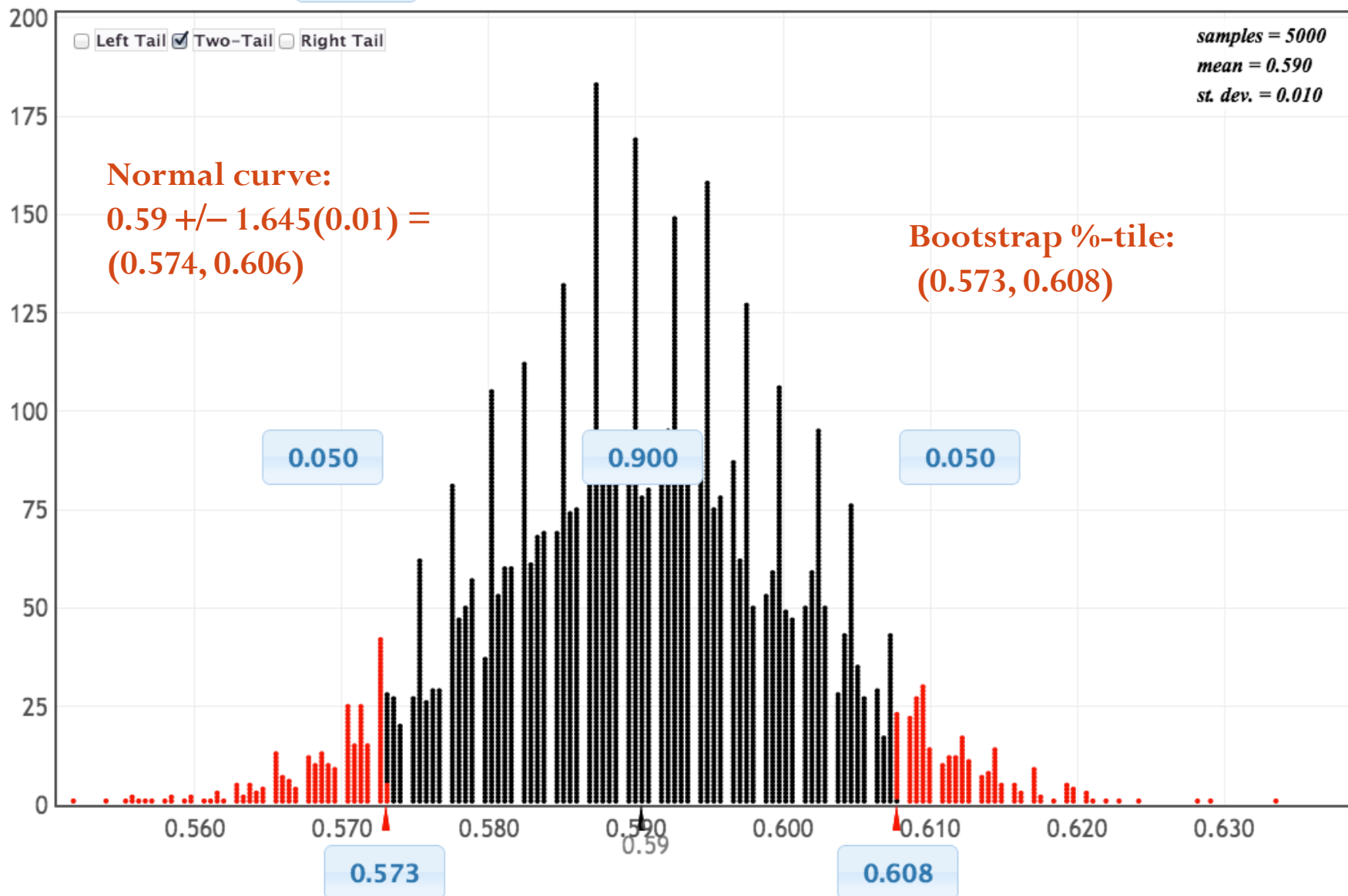
**$z^*=2.576$ for 99% confidence**

**$0.59 +/- 2.576(0.01) = (0.564, 0.616)$**

- Remember, more confidence = wider interval
- So how do these compare to the bootstrap CI?

# Global Warming: 90% CI



**Bootstrap Dotplot of** Proportion ▼
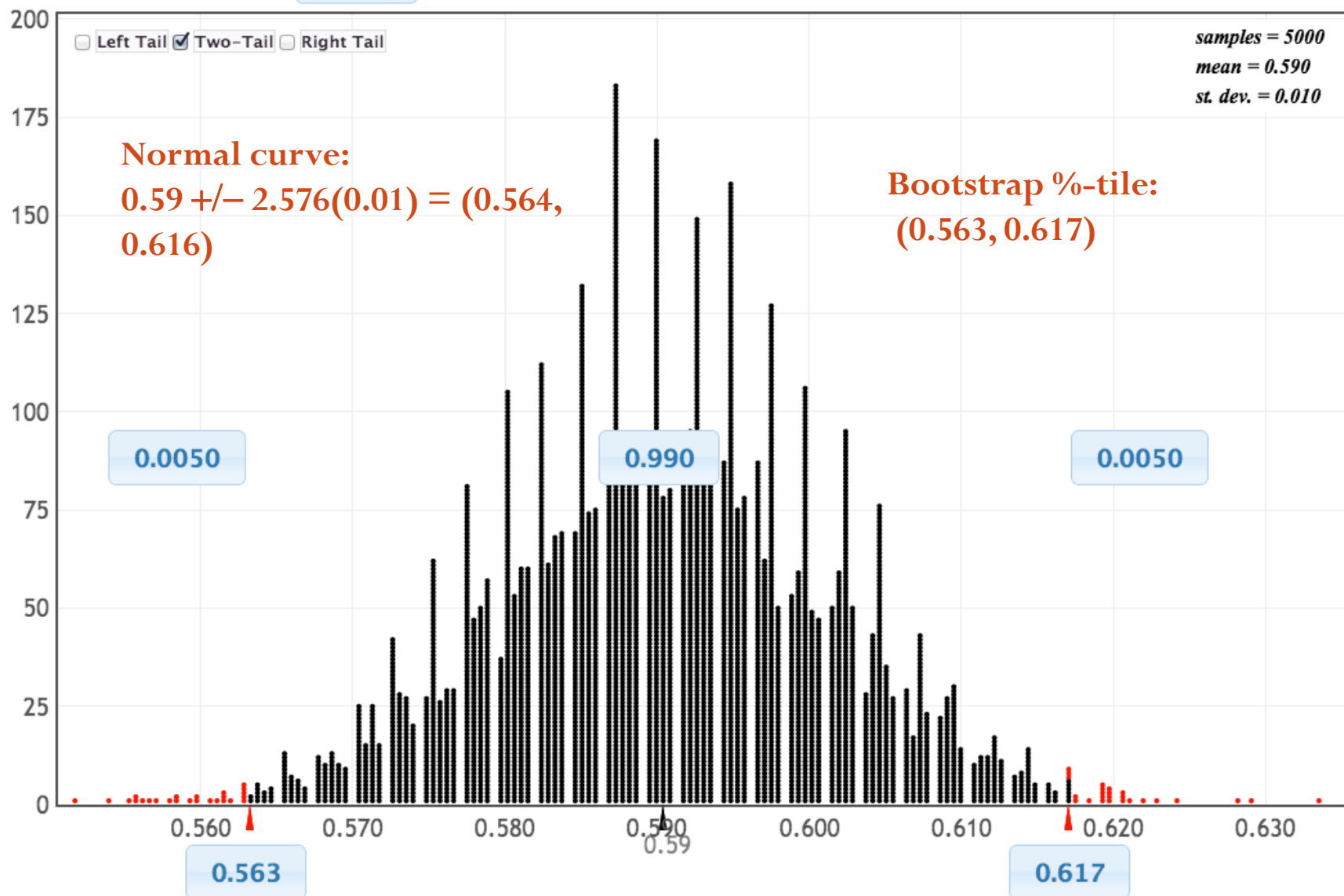
☐ Left Tail ☑ Two–Tail ☐ Right Tail

samples = 5000
mean = 0.590
st. dev. = 0.010

**Normal curve:**
**0.59 +/− 1.645(0.01) =**
**(0.574, 0.606)**

**Bootstrap %-tile:**
**(0.573, 0.608)**

0.050    0.900    0.050

0.573    0.608

# Global Warming: 99% CI

**Bootstrap Dotplot of** [Proportion ▾]



☐ Left Tail  ☑ Two-Tail  ☐ Right Tail

samples = 5000
mean = 0.590
st. dev. = 0.010

**Normal curve:**
**0.59 +/− 2.576(0.01) = (0.564, 0.616)**

**Bootstrap %-tile:**
**(0.563, 0.617)**

0.0050    0.990    0.0050

0.59

0.563    0.617

# Standard Error

- Wouldn't it be nice if we could compute the standard error *without* doing thousands of simulations?


- We can!!!

- Or rather, we'll be able to next class!