

# Inference for categorical data

Chunjie Nan

In August of 2012, news outlets ranging from the Washington Post to the Huffington Post ran a story about the rise of atheism in America. The source for the story was a poll that asked people, “Irrespective of whether you attend a place of worship or not, would you say you are a religious person, not a religious person or a convinced atheist?” This type of question, which asks people to classify themselves in one way or another, is common in polling and generates categorical data. In this lab we take a look at the atheism survey and explore what’s at play when making inference about population proportions using categorical data.

## The survey

To access the press release for the poll, conducted by WIN-Gallup International, click on the following link:

[https://github.com/jbryer/DATA606/blob/master/inst/labs/Lab6/more/Global\\_INDEX\\_of\\_Religiosity\\_and\\_Atheism\\_PR\\_\\_6.pdf](https://github.com/jbryer/DATA606/blob/master/inst/labs/Lab6/more/Global_INDEX_of_Religiosity_and_Atheism_PR__6.pdf)

Take a moment to review the report then address the following questions.

1. In the first paragraph, several key findings are reported. Do these percentages appear to be *sample statistics* (derived from the data sample) or *population parameters*?

Answer: The data is from a survey, which is sample proportion. so, the percentages are sample statistics.

2. The title of the report is “Global Index of Religiosity and Atheism”. To generalize the report’s findings to the global human population, what must we assume about the sampling method? Does that seem like a reasonable assumption?

Answer: yes, it is a reasonable assumption because the sample group is independent and it is randomly selected with enough observations.

## The data

Turn your attention to Table 6 (pages 15 and 16), which reports the sample size and response percentages for all 57 countries. While this is a useful format to summarize the data, we will base our analysis on the original data set of individual responses to the survey. Load this data set into R with the following command.

```
load("more/atheism.RData")
```

3. What does each row of Table 6 correspond to? What does each row of `atheism` correspond to?

Answer: It shows the results of the survey for each country. Each row of “atheism” corresponds to one observation.

To investigate the link between these two ways of organizing this data, take a look at the estimated proportion of atheists in the United States. Towards the bottom of Table 6, we see that this is 5%. We should be able to come to the same number using the `atheism` data.

4. Using the command below, create a new dataframe called `us12` that contains only the rows in `atheism` associated with respondents to the 2012 survey from the United States. Next, calculate the proportion of atheist responses. Does it agree with the percentage in Table 6? If not, why?

```
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(ggplot2)
us12 <- subset(atheism, nationality == "United States" & year == "2012")
prop_atheist <- us12 %>%
  summarise(atheist = mean(response == "atheist"))
prop_atheist

##   atheist
## 1 0.0499002
```

## Inference on proportions

As was hinted at in Exercise 1, Table 6 provides *statistics*, that is, calculations made from the sample of 51,927 people. What we'd like, though, is insight into the population *parameters*. You answer the question, “What proportion of people in your sample reported being atheists?” with a statistic; while the question “What proportion of people on earth would report being atheists” is answered with an estimate of the parameter.

The inferential tools for estimating population proportion are analogous to those used for means in the last chapter: the confidence interval and the hypothesis test.

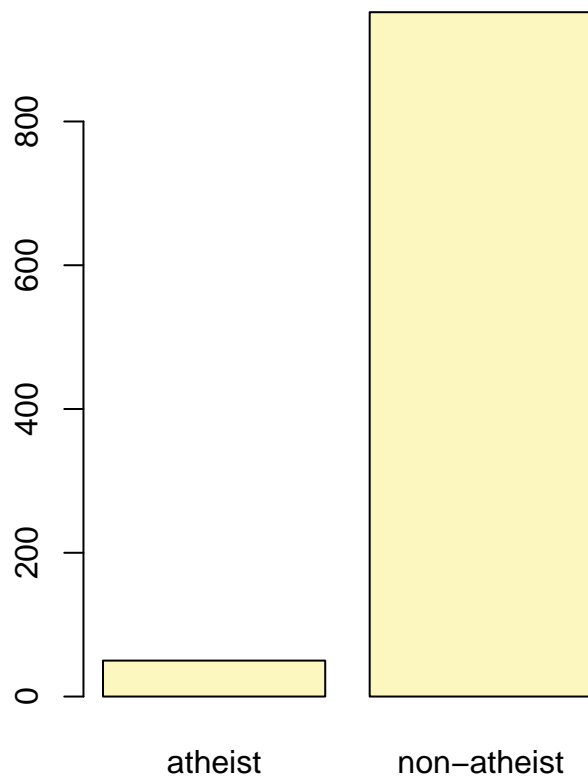
5. Write out the conditions for inference to construct a 95% confidence interval for the proportion of atheists in the United States in 2012. Are you confident all conditions are met?

Answer: yes, the sample is less than 10% of the population, and the sample observation is independent.

If the conditions for inference are reasonable, we can either calculate the standard error and construct the interval by hand, or allow the `inference` function to do it for us.

```
inference(us12$response, est = "proportion", type = "ci", method = "theoretical",
  success = "atheist")
```

```
## Single proportion -- success: atheist
## Summary statistics:
```



us12\$response

```
## p_hat = 0.0499 ; n = 1002
## Check conditions: number of successes = 50 ; number of failures = 952
## Standard error = 0.0069
## 95 % Confidence interval = ( 0.0364 , 0.0634 )
```

Note that since the goal is to construct an interval estimate for a proportion, it's necessary to specify what constitutes a “success”, which here is a response of "atheist".

Although formal confidence intervals and hypothesis tests don't show up in the report, suggestions of inference appear at the bottom of page 7: “In general, the error margin for surveys of this kind is  $\pm 3\text{-}5\%$  at 95% confidence”.

- Based on the R output, what is the margin of error for the estimate of the proportion of the proportion of atheists in US in 2012?

Answer:

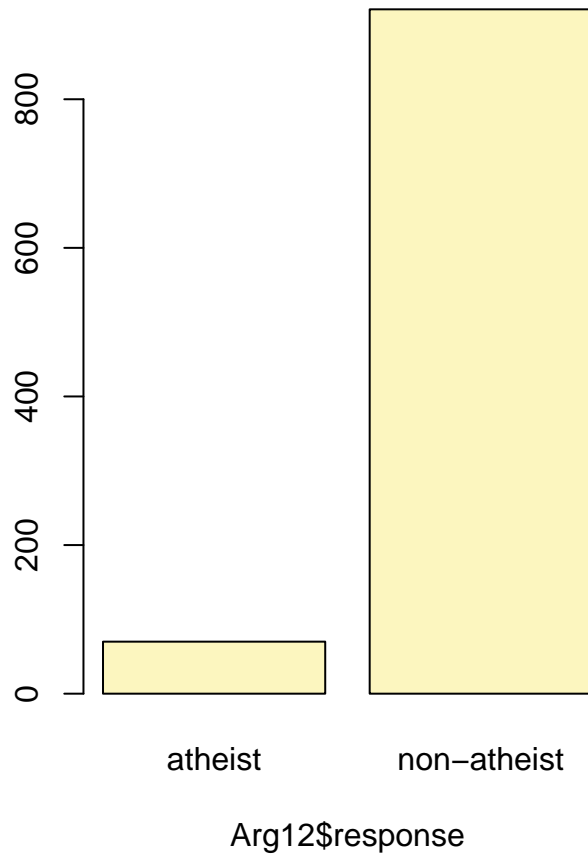
```
Margin_error<-1.96*0.0069
Margin_error
```

```
## [1] 0.013524
```

- Using the **inference** function, calculate confidence intervals for the proportion of atheists in 2012 in two other countries of your choice, and report the associated margins of error. Be sure to note whether the conditions for inference are met. It may be helpful to create new data sets for each of the two countries first, and then use these data sets in the **inference** function to construct the confidence intervals.

```
Arg12<- subset(atheism, nationality == "Argentina" & year == "2012")
inference(Arg12$response, est = "proportion" , type="ci", method = "theoretical", success="atheist")
```

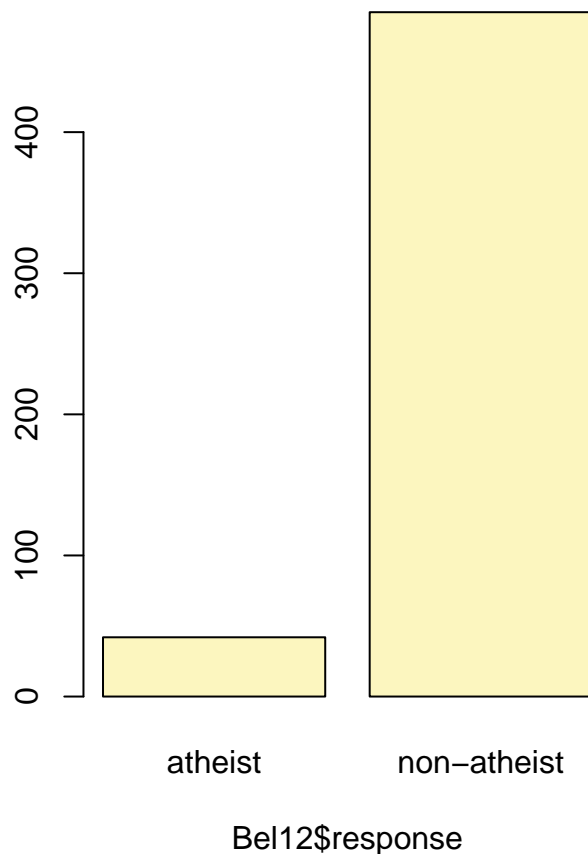
```
## Single proportion -- success: atheist
## Summary statistics:
```



```
## p_hat = 0.0706 ; n = 991
## Check conditions: number of successes = 70 ; number of failures = 921
## Standard error = 0.0081
## 95 % Confidence interval = ( 0.0547 , 0.0866 )
```

```
Bel12<- subset(atheism, nationality == "Belgium" & year == "2012")
inference(Bel12$response, est = "proportion" , type="ci", method = "theoretical", success="atheist")
```

```
## Single proportion -- success: atheist
## Summary statistics:
```



```
## p_hat = 0.0797 ; n = 527
## Check conditions: number of successes = 42 ; number of failures = 485
## Standard error = 0.0118
## 95 % Confidence interval = ( 0.0566 , 0.1028 )
```

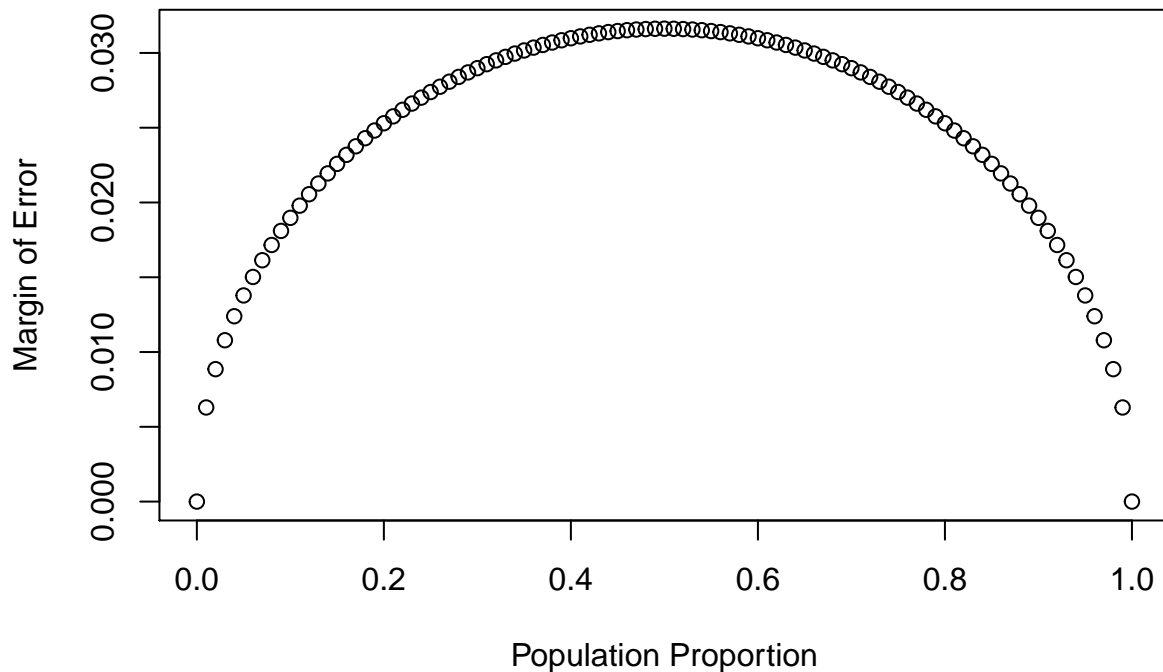
## How does the proportion affect the margin of error?

Imagine you've set out to survey 1000 people on two questions: are you female? and are you left-handed? Since both of these sample proportions were calculated from the same sample size, they should have the same margin of error, right? Wrong! While the margin of error does change with sample size, it is also affected by the proportion.

Think back to the formula for the standard error:  $SE = \sqrt{p(1-p)/n}$ . This is then used in the formula for the margin of error for a 95% confidence interval:  $ME = 1.96 \times SE = 1.96 \times \sqrt{p(1-p)/n}$ . Since the population proportion  $p$  is in this  $ME$  formula, it should make sense that the margin of error is in some way dependent on the population proportion. We can visualize this relationship by creating a plot of  $ME$  vs.  $p$ .

The first step is to make a vector  $\mathbf{p}$  that is a sequence from 0 to 1 with each number separated by 0.01. We can then create a vector of the margin of error ( $\mathbf{me}$ ) associated with each of these values of  $\mathbf{p}$  using the familiar approximate formula ( $ME = 2 \times SE$ ). Lastly, we plot the two vectors against each other to reveal their relationship.

```
n <- 1000
p <- seq(0, 1, 0.01)
me <- 2 * sqrt(p * (1 - p)/n)
plot(me ~ p, ylab = "Margin of Error", xlab = "Population Proportion")
```



8. Describe the relationship between  $p$  and  $me$ .

### Success-failure condition

The textbook emphasizes that you must always check conditions before making inference. For inference on proportions, the sample proportion can be assumed to be nearly normal if it is based upon a random sample of independent observations and if both  $np \geq 10$  and  $n(1 - p) \geq 10$ . This rule of thumb is easy enough to follow, but it makes one wonder: what's so special about the number 10?

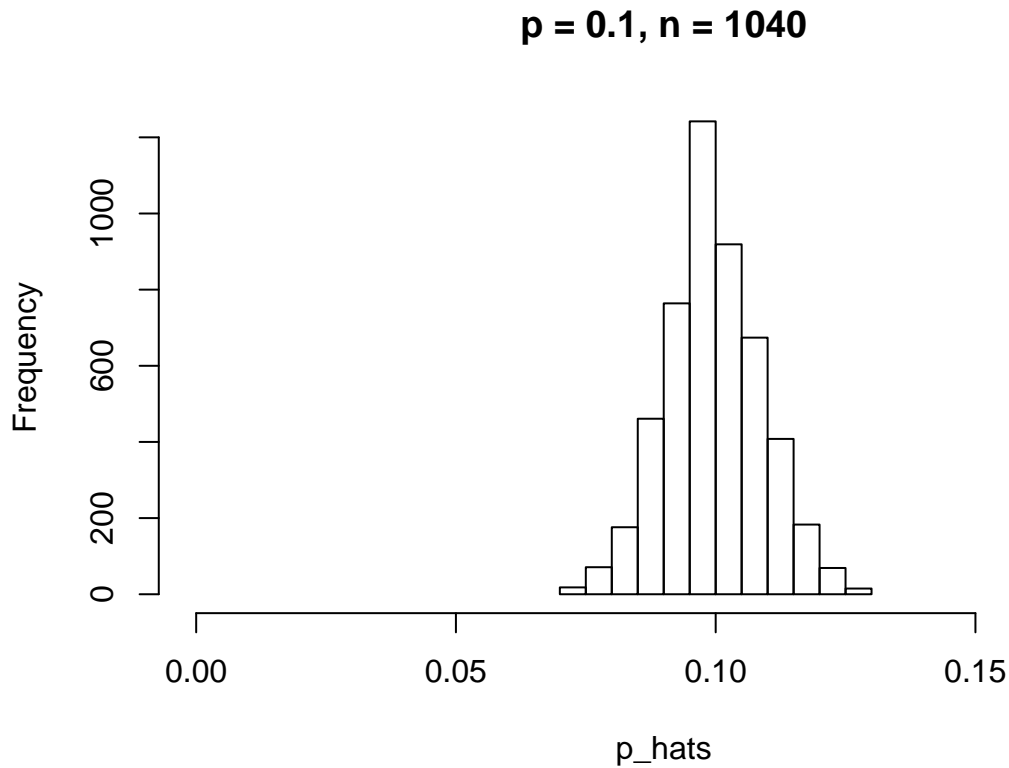
The short answer is: nothing. You could argue that we would be fine with 9 or that we really should be using 11. What is the “best” value for such a rule of thumb is, at least to some degree, arbitrary. However, when  $np$  and  $n(1 - p)$  reaches 10 the sampling distribution is sufficiently normal to use confidence intervals and hypothesis tests that are based on that approximation.

We can investigate the interplay between  $n$  and  $p$  and the shape of the sampling distribution by using simulations. To start off, we simulate the process of drawing 5000 samples of size 1040 from a population with a true atheist proportion of 0.1. For each of the 5000 samples we compute  $\hat{p}$  and then plot a histogram to visualize their distribution.

```
p <- 0.1
n <- 1040
p_hats <- rep(0, 5000)

for(i in 1:5000){
  samp <- sample(c("atheist", "non_atheist"), n, replace = TRUE, prob = c(p, 1-p))
  p_hats[i] <- sum(samp == "atheist")/n
}

hist(p_hats, main = "p = 0.1, n = 1040", xlim = c(0, 0.18))
```



These commands build up the sampling distribution of  $\hat{p}$  using the familiar **for** loop. You can read the sampling procedure for the first line of code inside the **for** loop as, “take a sample of size  $n$  with replacement from the choices of atheist and non-atheist with probabilities  $p$  and  $1 - p$ , respectively.” The second line in the loop says, “calculate the proportion of atheists in this sample and record this value.” The loop allows us to repeat this process 5,000 times to build a good representation of the sampling distribution.

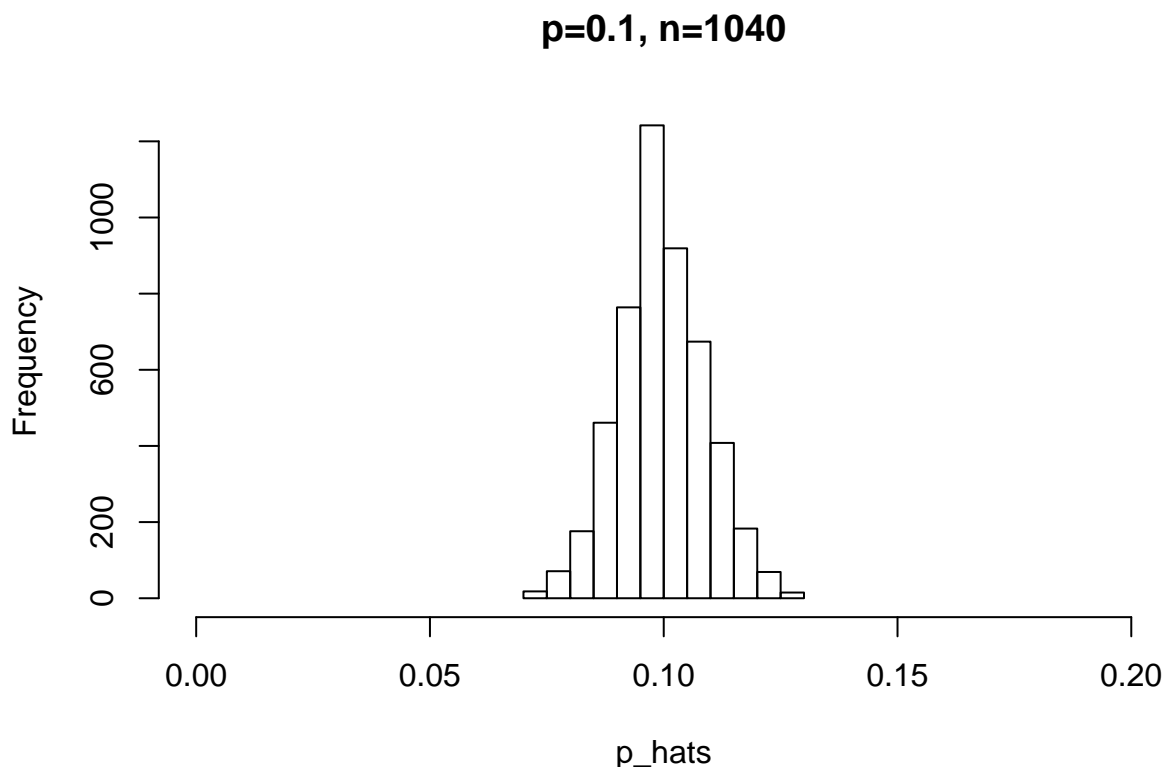
9. Describe the sampling distribution of sample proportions at  $n = 1040$  and  $p = 0.1$ . Be sure to note the center, spread, and shape.

*Hint:* Remember that R has functions such as **mean** to calculate summary statistics.

```
summary(p_hats)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.07019 0.09327 0.09904 0.09969 0.10577 0.12981
```

```
hist(p_hats, main= "p=0.1, n=1040",xlim = c(0,0.2))
```



10. Repeat the above simulation three more times but with modified sample sizes and proportions: for  $n = 400$  and  $p = 0.1$ ,  $n = 1040$  and  $p = 0.02$ , and  $n = 400$  and  $p = 0.02$ . Plot all four histograms together by running the `par(mfrow = c(2, 2))` command before creating the histograms. You may need to expand the plot window to accommodate the larger two-by-two plot. Describe the three new sampling distributions. Based on these limited plots, how does  $n$  appear to affect the distribution of  $\hat{p}$ ? How does  $p$  affect the sampling distribution?

Once you're done, you can reset the layout of the plotting window by using the command `par(mfrow = c(1, 1))` command or clicking on "Clear All" above the plotting window (if using RStudio). Note that the latter will get rid of all your previous plots.

```
p <- 0.1
n <- 400
p_hats1 <- rep(0, 5000)

for(i in 1:5000){
  samp1 <- sample(c("atheist", "non_atheist"), n, replace = TRUE, prob = c(p, 1-p))
  p_hats1[i] <- sum(samp1 == "atheist")/n
}
```

```
p <- 0.02
n <- 1040
p_hats2 <- rep(0, 5000)

for(i in 1:5000){
  samp2 <- sample(c("atheist", "non_atheist"), n, replace = TRUE, prob = c(p, 1-p))
  p_hats2[i] <- sum(samp2 == "atheist")/n
}
```



```

p <- 0.02
n <- 400
p_hats3 <- rep(0, 5000)

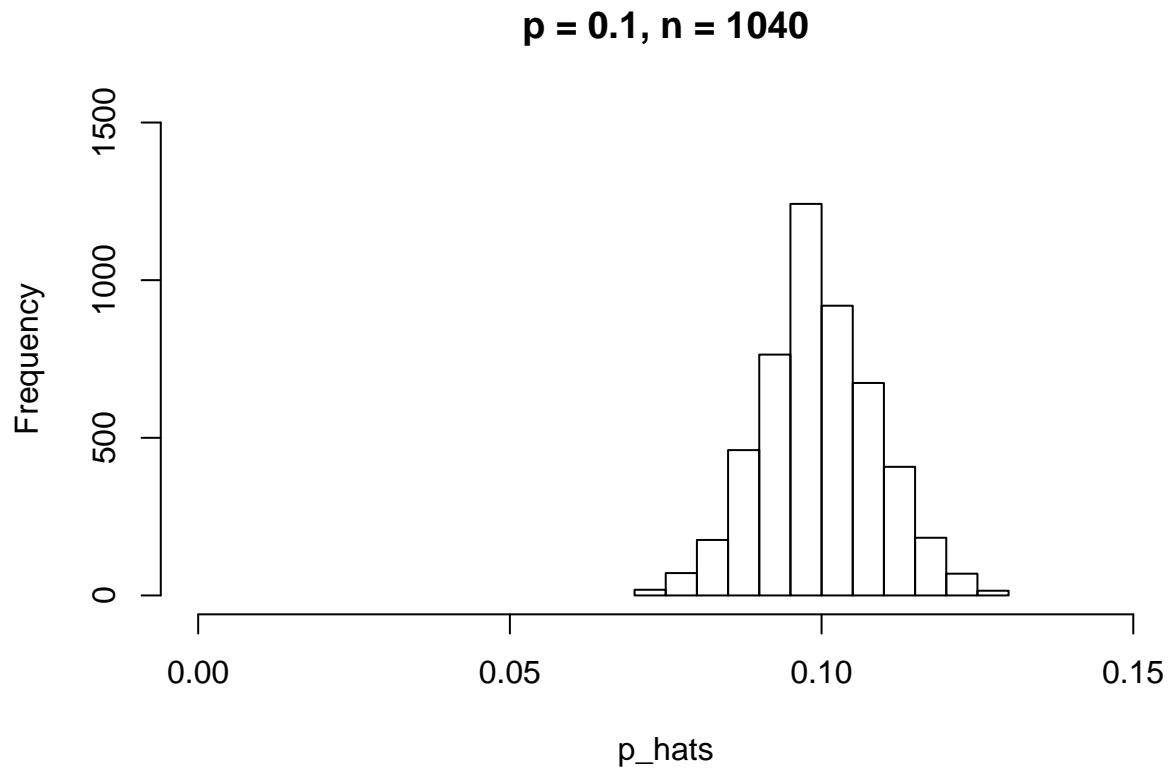
for(i in 1:5000){
  samp3 <- sample(c("atheist", "non_atheist"), n, replace = TRUE, prob = c(p, 1-p))
  p_hats3[i] <- sum(samp3 == "atheist")/n
}

```

```

hist(p_hats, main = "p = 0.1, n = 1040", xlim = c(0, 0.15), ylim = c(0, 1500))

```

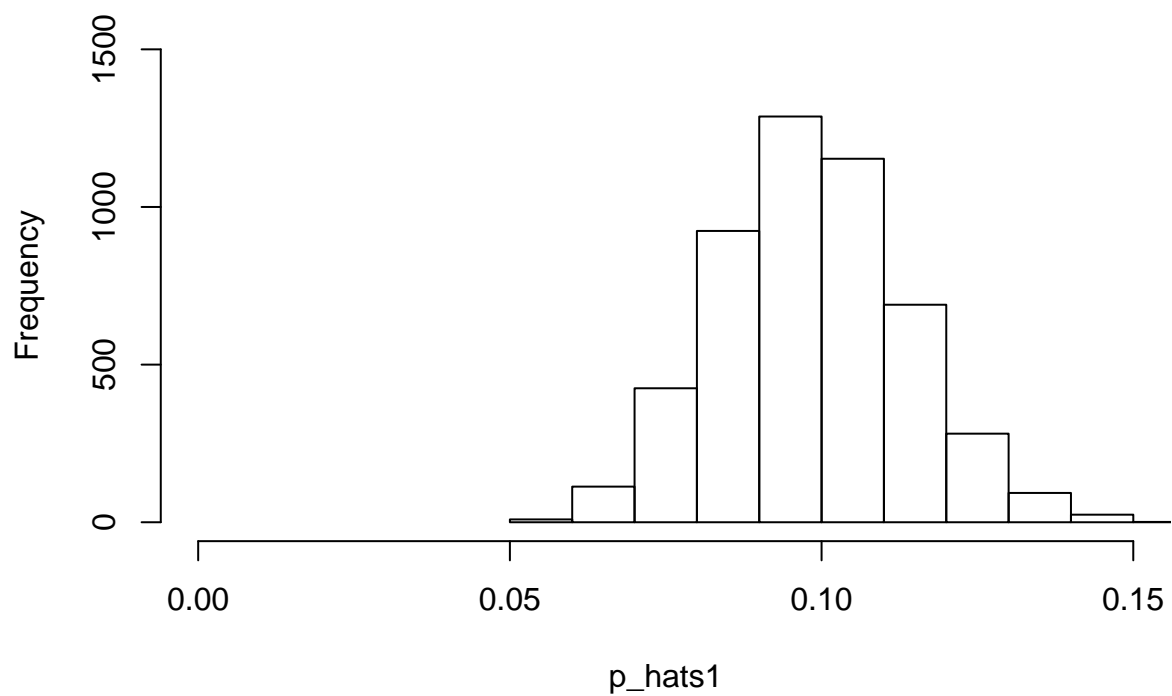


```

hist(p_hats1, main = "p = 0.1, n = 400", xlim = c(0, 0.15), ylim = c(0, 1500))

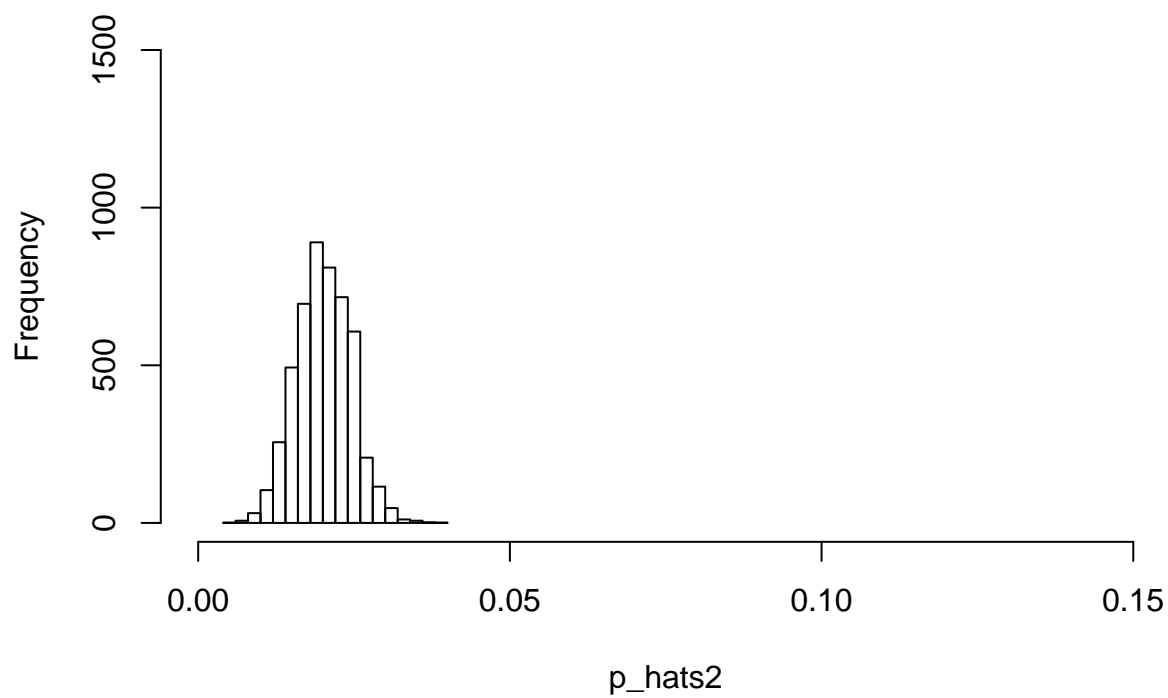
```

**p = 0.1, n = 400**

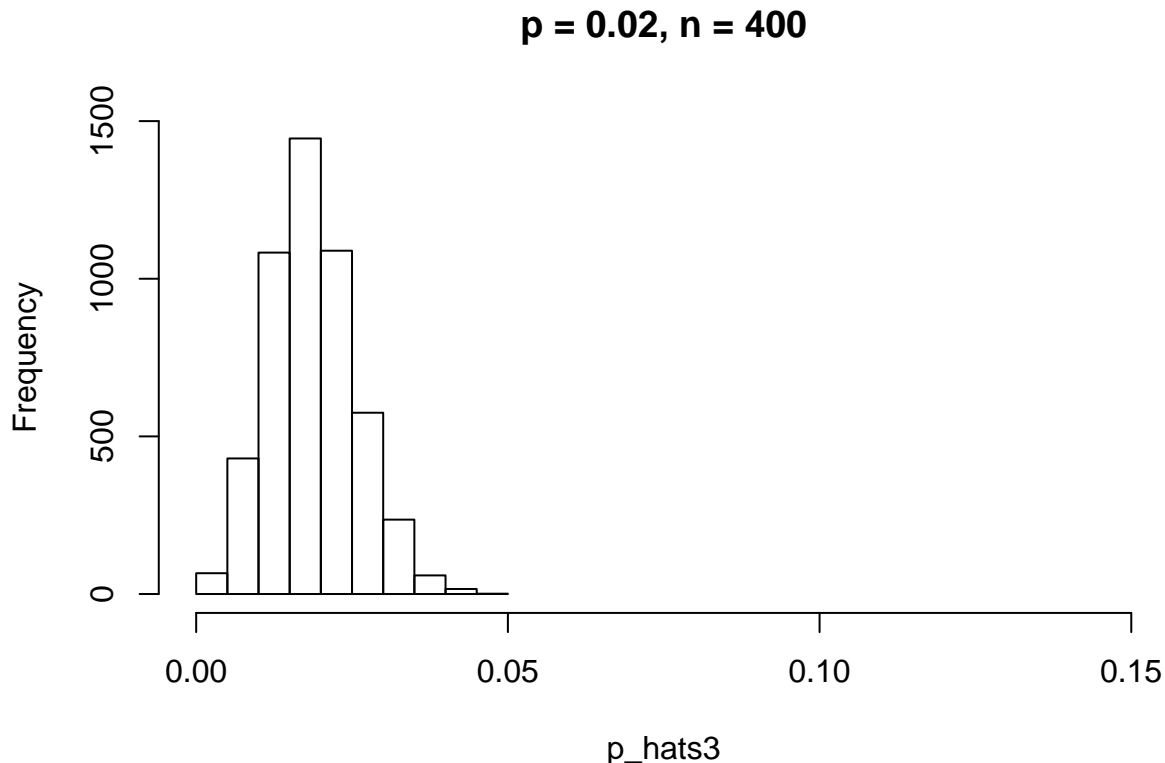


```
hist(p_hats2, main = "p = 0.02, n = 1040", xlim = c(0, 0.15), ylim = c(0, 1500))
```

**p = 0.02, n = 1040**



```
hist(p_hats3, main = "p = 0.02, n = 400", xlim = c(0, 0.15), ylim = c(0, 1500))
```



Answer: The larger  $n$ , the narrower distribution, and the larger  $p$  results the wider distribution.

11. If you refer to Table 6, you'll find that Australia has a sample proportion of 0.1 on a sample size of 1040, and that Ecuador has a sample proportion of 0.02 on 400 subjects. Let's suppose for this exercise that these point estimates are actually the truth. Then given the shape of their respective sampling distributions, do you think it is sensible to proceed with inference and report margin of errors, as the reports does?

Answer: The histograms have similar spread, so the margin of errors are also similar. \* \* \* ## On your own

The question of atheism was asked by WIN-Gallup International in a similar survey that was conducted in 2005. (We assume here that sample sizes have remained the same.) Table 4 on page 13 of the report summarizes survey results from 2005 and 2012 for 39 countries.

- Answer the following two questions using the **inference** function. As always, write out the hypotheses for any tests you conduct and outline the status of the conditions for inference.
  - a. Is there convincing evidence that Spain has seen a change in its atheism index between 2005 and 2012?  
*Hint:* Create a new data set for respondents from Spain. Form confidence intervals for the true proportion of athiests in both years, and determine whether they overlap.

Answer:

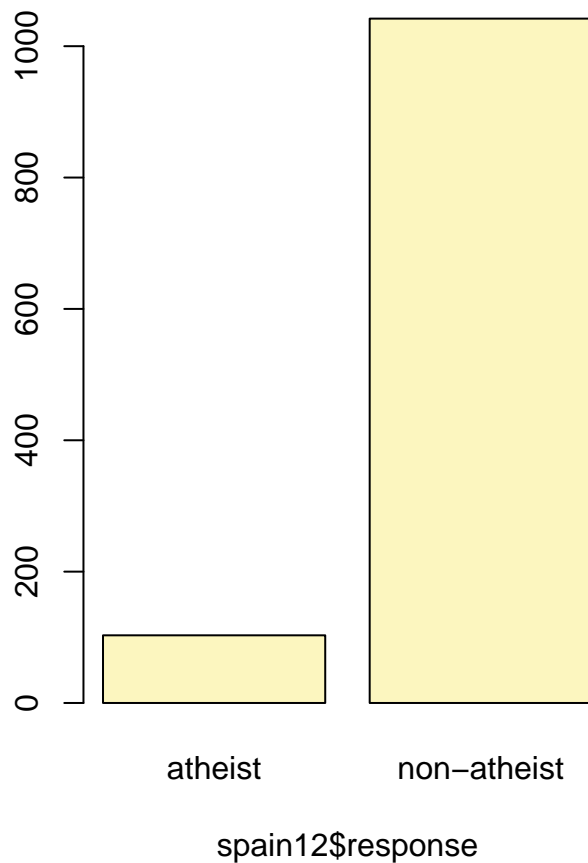
$H_0$ : There is no convincing evidence that Spain has seen a change in its atheism index between 2005 and 2012.

$H_A$ : There is convincing evidence that Spain has seen a change in its atheism index between 2005 and 2012.

```
spain12 <- subset(atheism, nationality == "Spain" & year == "2012")
spain05 <- subset(atheism, nationality == "Spain" & year == "2005")

inference(spain12$response, est = "proportion", type = "ci", method = "theoretical",
          success = "atheist")
```

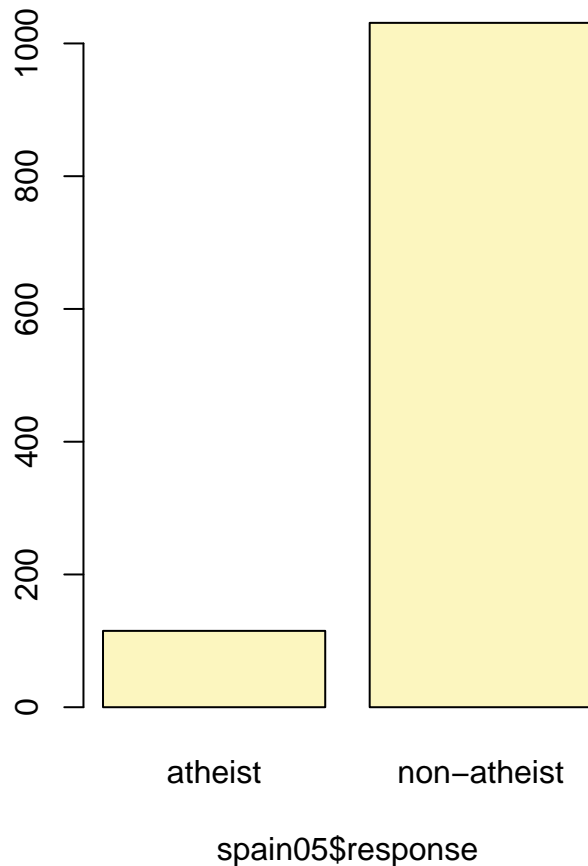
```
## Single proportion -- success: atheist
## Summary statistics:
```



```
## p_hat = 0.09 ; n = 1145
## Check conditions: number of successes = 103 ; number of failures = 1042
## Standard error = 0.0085
## 95 % Confidence interval = ( 0.0734 , 0.1065 )
```

```
inference(spain05$response, est = "proportion", type = "ci", method = "theoretical",
          success = "atheist")
```

```
## Single proportion -- success: atheist
## Summary statistics:
```



```
## p_hat = 0.1003 ; n = 1146
## Check conditions: number of successes = 115 ; number of failures = 1031
## Standard error = 0.0089
## 95 % Confidence interval = ( 0.083 , 0.1177 )
```

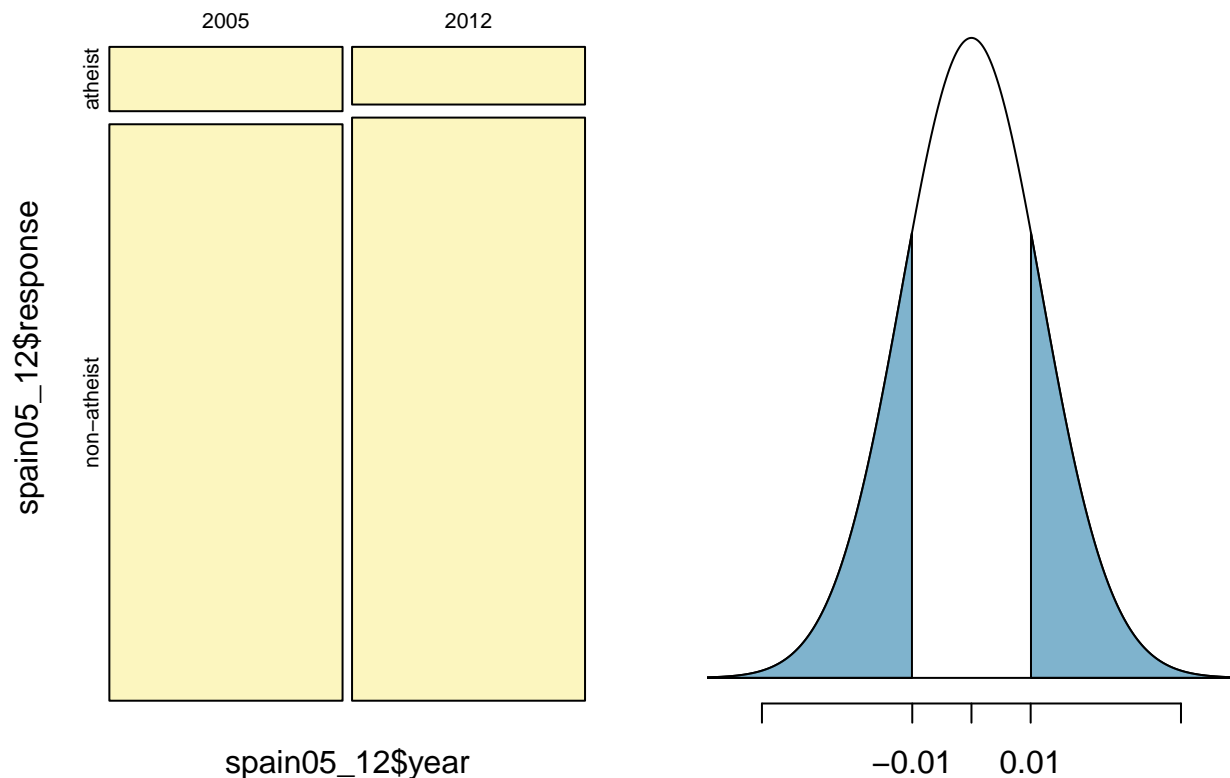
Answer: The 95% of confident intervals are close to each other.

```
spain05_12 <- subset(atheism, nationality == "Spain" & year == "2005" | nationality == "Spain" & year == "2012")
inference(y = spain05_12$response, x = spain05_12$year, est = "proportion", type = "ht", null = 0, alternative = "two.sided")
```

```
## Warning: Explanatory variable was numerical, it has been converted to
## categorical. In order to avoid this warning, first convert your explanatory
## variable to a categorical variable using the as.factor() function.
```

```
## Response variable: categorical, Explanatory variable: categorical
## Two categorical variables
## Difference between two proportions -- success: atheist
## Summary statistics:
##           x
## y          2005 2012 Sum
## atheist        115  103 218
## non-atheist    1031 1042 2073
## Sum            1146 1145 2291
```

```
## Observed difference between proportions (2005-2012) = 0.0104
##
## H0: p_2005 - p_2012 = 0
## HA: p_2005 - p_2012 != 0
## Pooled proportion = 0.0952
## Check conditions:
##   2005 : number of expected successes = 109 ; number of expected failures = 1037
##   2012 : number of expected successes = 109 ; number of expected failures = 1036
## Standard error = 0.012
## Test statistic: Z = 0.848
## p-value = 0.3966
```



Answer: The p-value is 0.3699 which is beyond 0.05, so we fail to reject the H<sub>0</sub>. so, Spain is less likely to change in its atheism between 2005 and 2012.

**\*\*b.\*\*** Is there convincing evidence that the United States has seen a change in its atheism index between 2005 and 2012?

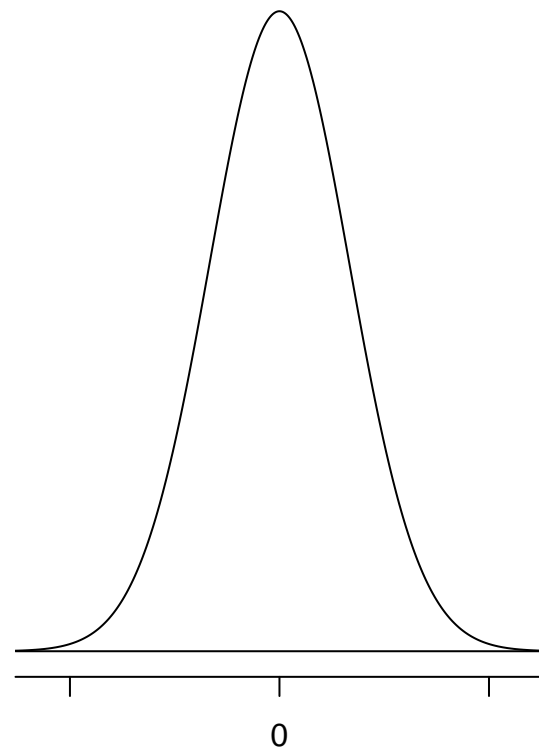
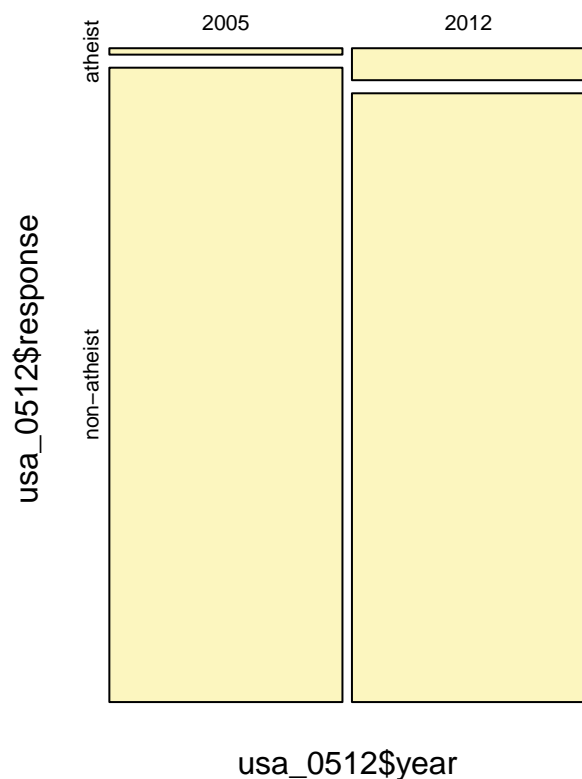
```
usa_0512 <- subset(atheism, nationality == "United States" & year == "2005" | nationality == "United States" & year == "2012")
inference(y = usa_0512$response, x = usa_0512$year, est = "proportion", type = "ht", null = 0, alternative = "two.sided")
```

```
## Warning: Explanatory variable was numerical, it has been converted to
## categorical. In order to avoid this warning, first convert your explanatory
## variable to a categorical variable using the as.factor() function.
```

```
## Response variable: categorical, Explanatory variable: categorical
## Two categorical variables
```

```
## Difference between two proportions -- success: atheist
## Summary statistics:
##           x
## y      2005 2012  Sum
## atheist      10   50   60
## non-atheist  992  952 1944
## Sum        1002 1002 2004

## Observed difference between proportions (2005-2012) = -0.0399
##
## H0: p_2005 - p_2012 = 0
## HA: p_2005 - p_2012 != 0
## Pooled proportion = 0.0299
## Check conditions:
##   2005 : number of expected successes = 30 ; number of expected failures = 972
##   2012 : number of expected successes = 30 ; number of expected failures = 972
## Standard error = 0.008
## Test statistic: Z = -5.243
## p-value = 0
```



Answer: The p-value is less than 0.05, so we reject the H<sub>0</sub>. so, USA is likely change in its atheism between 2005 and 2012.

- If in fact there has been no change in the atheism index in the countries listed in Table 4, in how many of those countries would you expect to detect a change (at a significance level of 0.05) simply by chance?

*Hint:* Look in the textbook index under Type 1 error.

Answer: A type 1 error is rejecting the null hypothesis when  $H_0$  is actually true. We do not want to incorrectly reject  $H_0$  more than 5% of the time. There are 39 countries listed in Table 4, we can count 5% of 39. So, 2 countries would be expected to detect a change in atheism just by chance.

- Suppose you're hired by the local government to estimate the proportion of residents that attend a religious service on a weekly basis. According to the guidelines, the estimate must have a margin of error no greater than 1% with 95% confidence. You have no idea what to expect for  $p$ . How many people would you have to sample to ensure that you are within the guidelines?

*Hint:* Refer to your plot of the relationship between  $p$  and margin of error. Do not use the data set to answer this question.

Answer: We choose the largest p-value and marginal error to determine the situation.

```
p<-0.5
z<-1.96
marginal_error<-0.01

n<- (p*(1-p)*z^2)/marginal_error^2
n
```

```
## [1] 9604
```

Answer: Therefore, we need at least 9604 people included in the sample.