

Inference for categorical data

Ilya Kats

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:

http://www.wingia.com/web/files/richeditor/filemanager/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*?

The percentages are generated from results of the survey which is a sample of the population, so they 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?

We must assume that the observations are independent. It is a reasonable assumption provided that the survey picked individuals at random and at 51,927 observations we can confidently say that the sample is less than 10% of the population. Since we are comparing countries and religiosity there may be minimum sample size requirements for inference for some smaller countries with small religious minorities.

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?

Each row in table 6 corresponds to results of the survey for each country. Each row of `atheism` corresponds to one observation (one individual surveyed).

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.

- 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?

```
us12 <- subset(atheism, nationality == "United States" & year == "2012")

# Remove unused levels
us12$nationality <- as.factor(as.character(us12$nationality))

# Get proportions
( us12prop <- prop.table(table(us12$nationality, us12$response)) )
```

```
##
##               atheist non-atheist
##   United States 0.0499002   0.9500998
```

As calculated, the proportion of atheists in the United States is 0.0499002 which slightly differs from the report's value of 0.05 because of rounding.

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.

- 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?

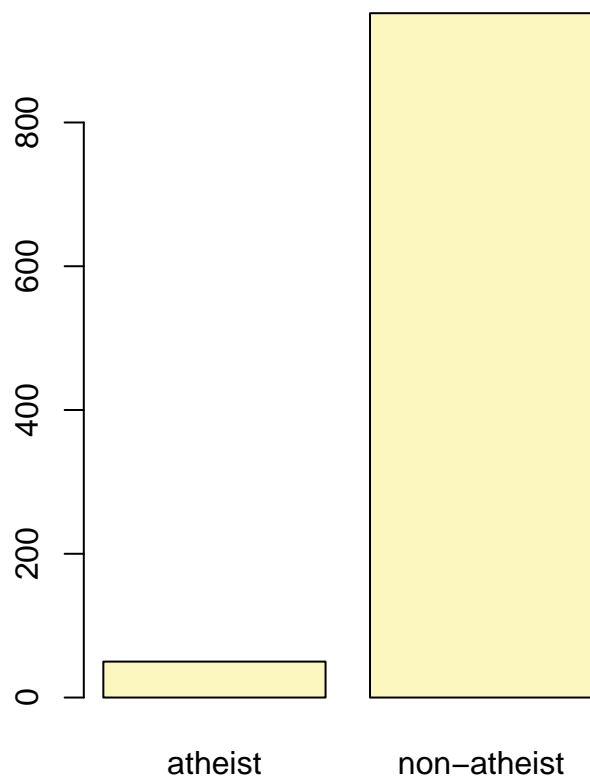
Observations must be independent. Assuming that individuals were selected using a simple random sample and considering that the sample is less than 10% of the population, this condition is satisfied.

Observations must come from a nearly normal distribution. Considering percentage of atheists at 0.05 and number of observations at 1,002, the observed number of atheists is 50 which is greater than 10. Assumption of nearly normal distribution is reasonable.

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?

Confidence interval is (0.0364, 0.0634), so margin of error is $(0.0634 - 0.0364) / 2 = 0.0135$. Alternatively, margin of error is $1.96 * \text{Standard Error} = 1.96 * 0.0069 = 0.013524$.

- 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.

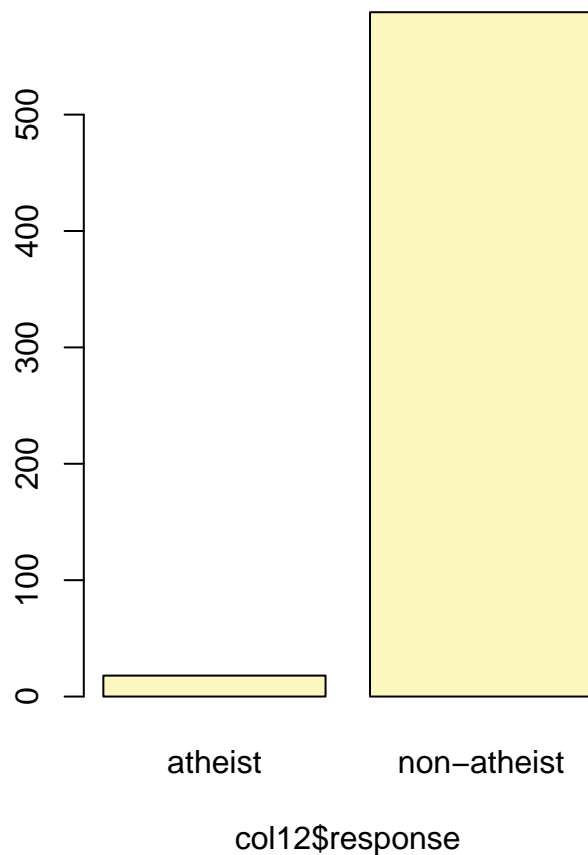
```
# Generate proportions for Colombia
col12 <- subset(atheism, nationality == "Colombia" & year == "2012")
```

```
col12$nationality <- as.factor(as.character(col12$nationality))
( col12prop <- prop.table(table(col12$nationality, col12$response)) )
```

```
##
##               atheist non-atheist
##   Colombia 0.02970297 0.97029703
```

```
inference(col12$response, est = "proportion", type = "ci", method = "theoretical",
          success = "atheist")
```

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



```
## p_hat = 0.0297 ; n = 606
## Check conditions: number of successes = 18 ; number of failures = 588
## Standard error = 0.0069
## 95 % Confidence interval = ( 0.0162 , 0.0432 )
```

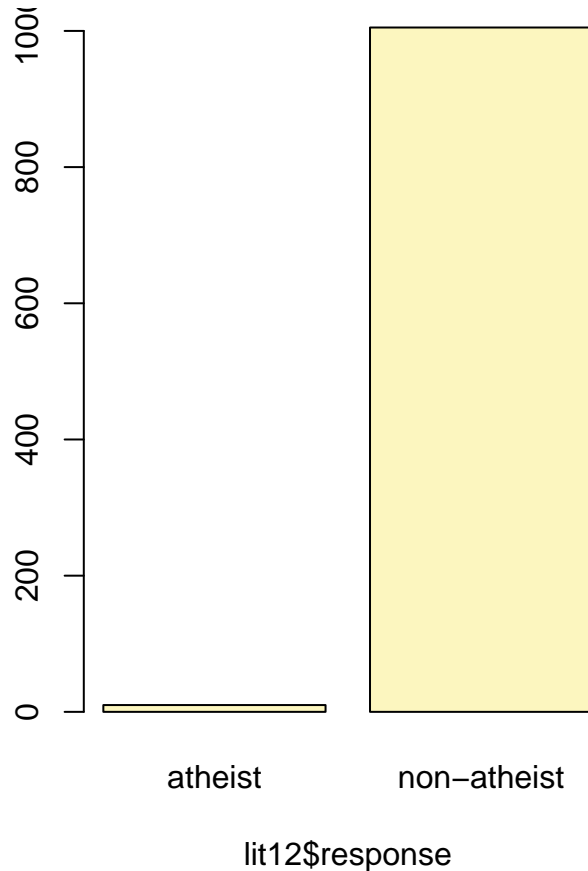
Margin of error for Colombia is $(0.0432 - 0.0162) / 2 = 0.0135$.

```
# Generate proportions for Lithuania
lit12 <- subset(atheism, nationality == "Lithuania" & year == "2012")
lit12$nationality <- as.factor(as.character(lit12$nationality))
( lit12prop <- prop.table(table(lit12$nationality, lit12$response)) )
```

```
##
##               atheist non-atheist
##   Lithuania 0.009852217 0.990147783
```

```
inference(lit12$response, est = "proportion", type = "ci", method = "theoretical",
          success = "atheist")
```

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



```
## p_hat = 0.0099 ; n = 1015
## Check conditions: number of successes = 10 ; number of failures = 1005
## Standard error = 0.0031
## 95 % Confidence interval = ( 0.0038 , 0.0159 )
```

Margin of error for Lithuania is $(0.0159 - 0.0038) / 2 = 0.00605$.

Similar to conditions for US sample, we can assume that samples for Colombia and Lithuania are independent and follow nearly normal distribution. Although the percentage of atheists in Lithuania is small, the number of atheists in the sample is 10, so it is borderline acceptable to assume nearly normal distribution.

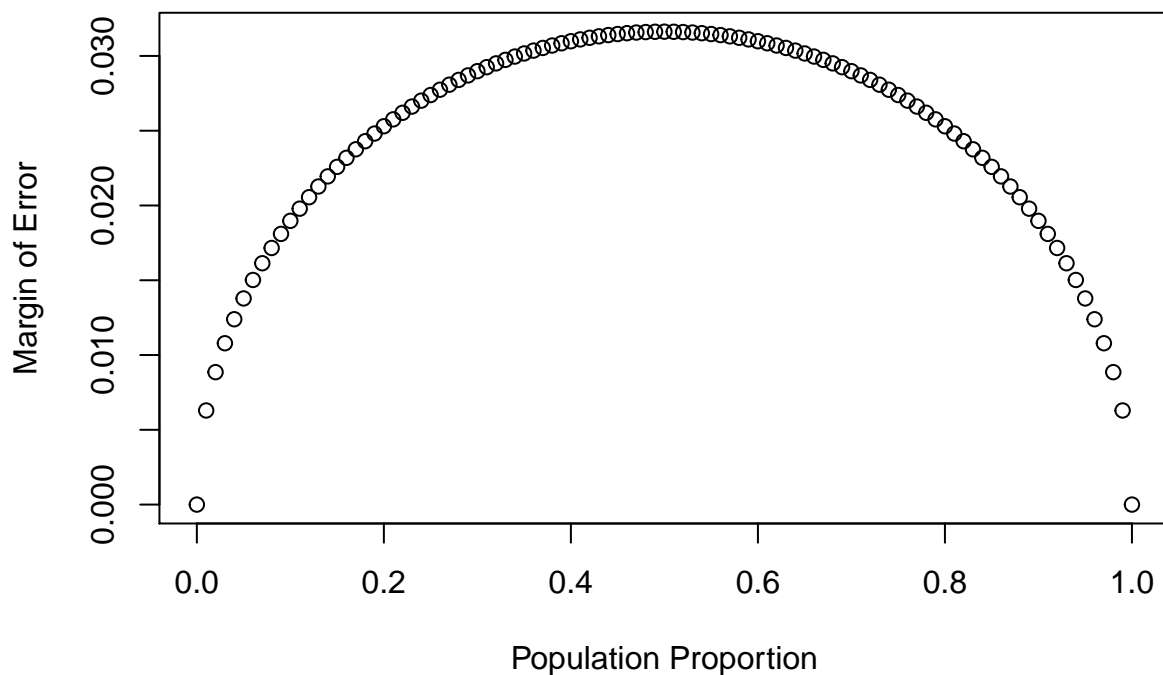
How does the proportion affect the margin of error?

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 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 (me) associated with each of these values of 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 .

The relationship is quadratic. As population proportion increases from 0 to 0.5, the margin of error also increases with the highest value at $p = 0.5$. As population proportion continues to increase past 0.5 the reverse is true and the margin of error decreases down to 0 at $p = 1$. Because margin of error is based on $p * (1 - p)$, the two halves of the relationship - with p from 0 to 0.5 and from 0.5 to 1 - are mirror images of each other.

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?

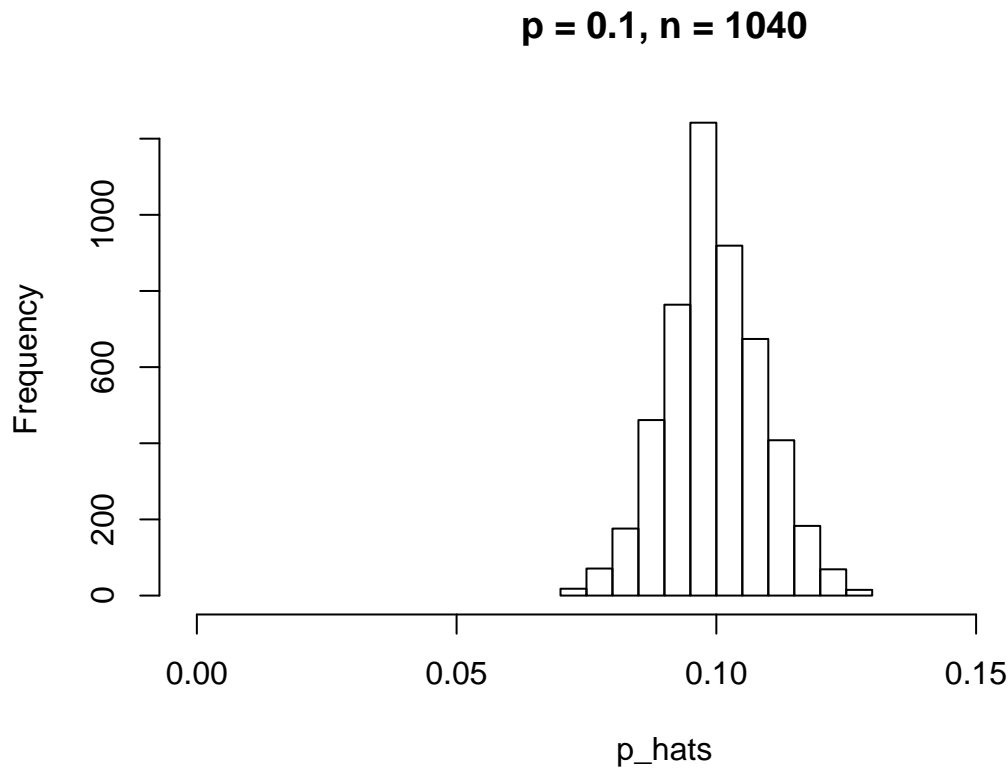
The short answer is: nothing. You could argue that we would be fine with 9 or that we really should be using 11. What is the “best” value for such a rule of thumb is, at least to some degree, arbitrary. However, when np and $n(1 - p)$ reaches 10 the sampling distribution is sufficiently normal to use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



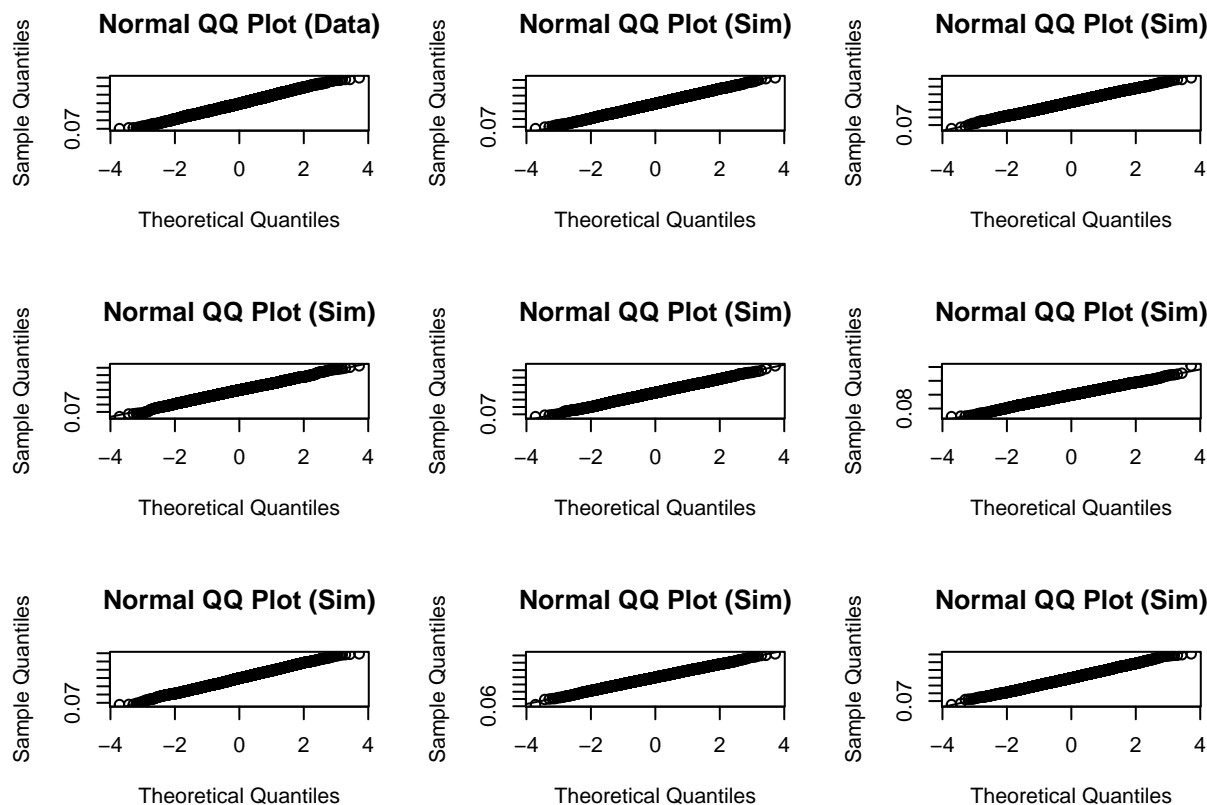
These commands build up the sampling distribution of \hat{p} using the familiar `for` loop. You can read the sampling procedure for the first line of code inside the `for` loop as, “take a sample of size n with replacement from the choices of atheist and non-atheist with probabilities p and $1 - p$, 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.

The sampling distribution is unimodal and symmetrical. It is centered around mean = 0.09969 with standard deviation of 0.0092874. Based on QQ-plots it appears to be nearly normal.

```
library("DATA606")
qqnormsim(p_hats)
```



-
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?

```
p_too <- c(0.1, 0.1, 0.02, 0.02)
n_too <- c(1040, 400, 1040, 400)
```



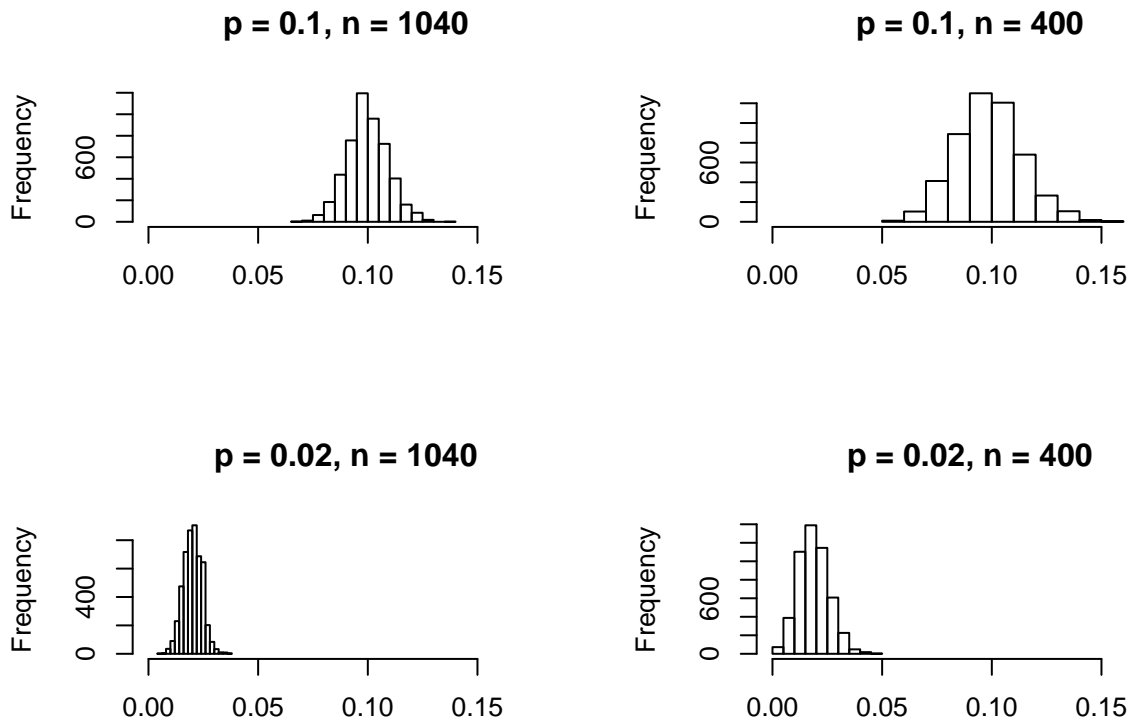
```

p_hats_too <- data.frame(c(rep(0, 5000)), c(rep(0, 5000)), c(rep(0, 5000)), c(rep(0, 5000)))

for(i in 1:5000){
  samp <- sample(c("atheist", "non_atheist"), n_too[1], replace = TRUE, prob = c(p_too[1], 1-p_too[1]))
  p_hats_too[i, 1] <- sum(samp == "atheist")/n_too[1]
  samp <- sample(c("atheist", "non_atheist"), n_too[2], replace = TRUE, prob = c(p_too[2], 1-p_too[2]))
  p_hats_too[i, 2] <- sum(samp == "atheist")/n_too[2]
  samp <- sample(c("atheist", "non_atheist"), n_too[3], replace = TRUE, prob = c(p_too[3], 1-p_too[3]))
  p_hats_too[i, 3] <- sum(samp == "atheist")/n_too[3]
  samp <- sample(c("atheist", "non_atheist"), n_too[4], replace = TRUE, prob = c(p_too[4], 1-p_too[4]))
  p_hats_too[i, 4] <- sum(samp == "atheist")/n_too[4]
}

par(mfrow = c(2, 2))
hist(p_hats_too[, 1], xlab = "", main = "p = 0.1, n = 1040", xlim = c(0, 0.18))
hist(p_hats_too[, 2], xlab = "", main = "p = 0.1, n = 400", xlim = c(0, 0.18))
hist(p_hats_too[, 3], xlab = "", main = "p = 0.02, n = 1040", xlim = c(0, 0.18))
hist(p_hats_too[, 4], xlab = "", main = "p = 0.02, n = 400", xlim = c(0, 0.18))

```



```

par(mfrow = c(1, 1))

```

All distributions are unimodal and symmetrical. The last one appears to have some skew, but all others appear to be nearly normal. Based on histograms above, n affects the spread of the distribution and p affects the center. Higher n value means the distribution is more narrow. And distributions are centered around p .

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.

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?

For Ecuador the sample proportion of atheists is 0.02 with 400 observations, so the number of atheists in the sample is 8, which is not enough to assume nearly normal distribution and it may not be sensible to proceed with inference. For Australia, with proportion of 0.1 and sample size of 1,040, the number of atheists in the sample is 104, so it is sensible to assume nearly normal distribution and proceed with inference.

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 atheists in both years, and determine whether they overlap.

```
spain05 <- subset(atheism, nationality == "Spain" & year == "2005")
spain05$nationality <- as.factor(as.character(spain05$nationality))
table(spain05$nationality, spain05$response)
```

```
##
##          atheist non-atheist
##   Spain      115      1031
```

```
spain12 <- subset(atheism, nationality == "Spain" & year == "2012")
spain12$nationality <- as.factor(as.character(spain12$nationality))
table(spain12$nationality, spain12$response)
```

```
##
##          atheist non-atheist
##   Spain      103      1042
```

As we established earlier in the lab, we can assume observations to be independent. The number of atheists in 2005 is 115 and in 2012 it is 103. Both are greater than 10, so we can assume near normal distribution.

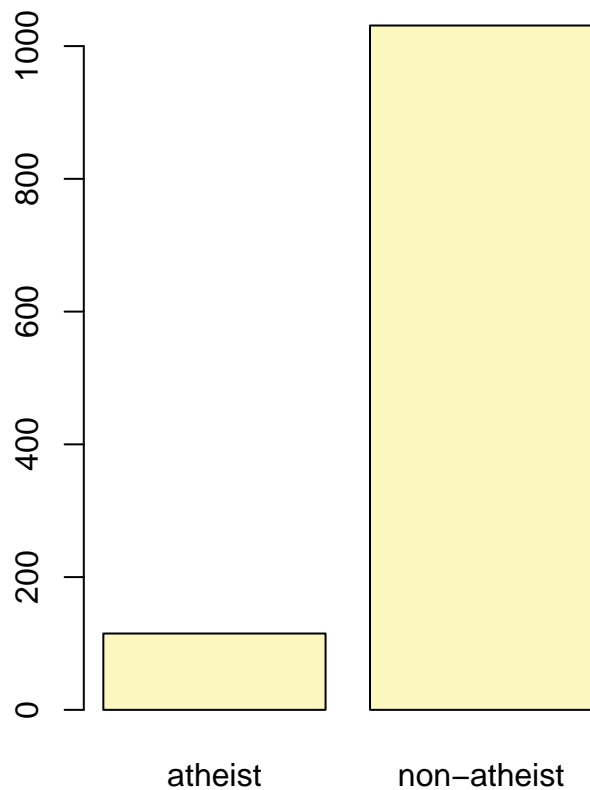
H_0 : The number of atheists in Spain did not change between 2005 and 2012, or $p_{12} = p_{05} = 0.1$.

H_A : The number of atheists in Spain changed between 2005 and 2012, or $p_{12} \neq 0.1$.

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

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

```
## Summary statistics:
```



spain05\$response

```
## p_hat = 0.1003 ; n = 1146
```

```
## Check conditions: number of successes = 115 ; number of failures = 1031
```

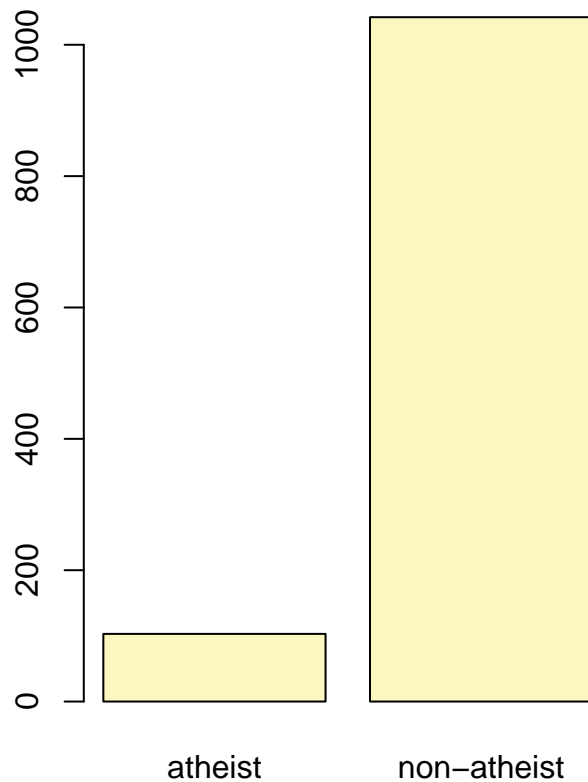
```
## Standard error = 0.0089
```

```
## 95 % Confidence interval = ( 0.083 , 0.1177 )
```

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

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

```
## Summary statistics:
```



spain12\$response

```
## 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 )
```

There is significant overlap between confidence interval for 2005 sample and 2012 sample. Additionally, $p_{12} = 0.09$ and it is within the confidence interval for 2005 - (0.081, 0.1177) - so we fail to reject the null hypothesis. The change in atheism is likely due to chance.

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

```
us05 <- subset(atheism, nationality == "United States" & year == "2005")
us05$nationality <- as.factor(as.character(us05$nationality))
table(us05$nationality, us05$response)
```

```
##
##           atheist non-atheist
## United States      10       992
```

```
table(us12$nationality, us12$response)
```

```
##
##           atheist non-atheist
## United States      50       952
```

As we established earlier in the lab, we can assume observations to be independent. The number of atheists in 2005 is 10 and in 2012 it is 50. The number of atheists in 2005 is borderline enough to assume near normal distribution.

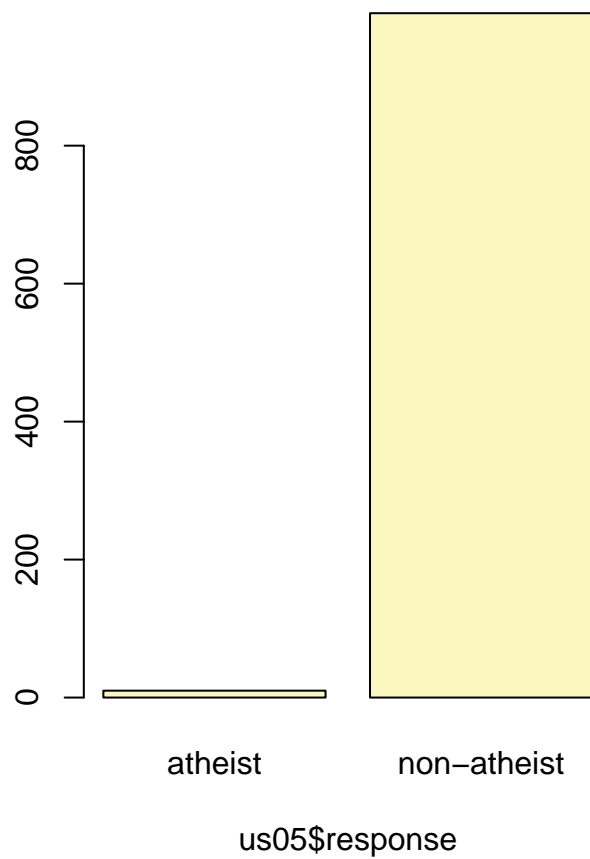
H_0 : The number of atheists in the United States did not change between 2005 and 2012, or $p_{12} = p_{05} = 0.01$.

H_A : The number of atheists in the United States changed between 2005 and 2012, or $p_{12} \neq 0.01$.

```
inference(us05$response, est = "proportion", type = "ci", method = "theoretical",
          success = "atheist")
```

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

```
## Summary statistics:
```



```
## p_hat = 0.01 ; n = 1002
```

```
## Check conditions: number of successes = 10 ; number of failures = 992
```

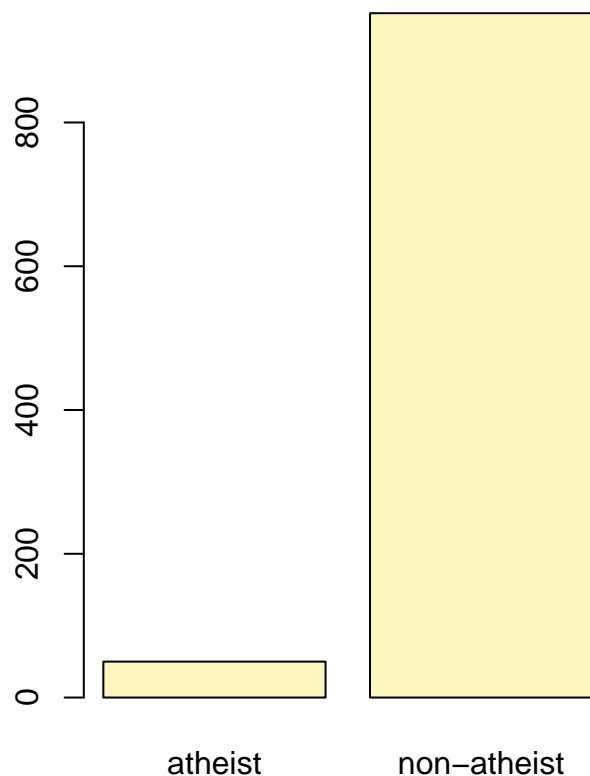
```
## Standard error = 0.0031
```

```
## 95 % Confidence interval = ( 0.0038 , 0.0161 )
```

```
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 )
```

There is no overlap between confidence interval for 2005 sample and 2012 sample. Additionally, $p_{12} = 0.05$ and it is outside of the confidence interval for 2005 - (0.0038, 0.0161) - so we reject the null hypothesis. The change in atheism is not likely due to chance.

- 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.

If there has been no change in the atheism index, but we detect a change due to chance and reject a null hypothesis even though it is true, that means we have made a Type I error. At a significant level of 0.05 and considering that we have 39 countries, we would expect to make a Type I error with $39 * 0.05 = 1.95$ or, rounding up, with 2 countries.

- 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.

If the margin of error is 0.01, then at 95% confidence $SE = \frac{0.01}{1.96} = 0.0051$. Since we do not know the value of p , we assume the worst case scenario with $p = 0.5$. Considering, $SE = \sqrt{\frac{p(1-p)}{n}}$, then $n = \frac{p(1-p)}{SE^2} = \frac{0.5*0.5}{0.0051^2} = 9604$. Sample must include at least 9,604 people. This is a worst

case scenario. It is possible that there is some good estimate of what proportion of residents attend services, so p can be lowered in the above calculation.

This is a product of OpenIntro that is released under a Creative Commons Attribution-ShareAlike 3.0 Unported. This lab was written for OpenIntro by Andrew Bray and Mine Çetinkaya-Rundel.