Lab 5: Inference for numerical data

North Carolina births

In 2004, the state of North Carolina released a large data set containing information on births recorded in this state. This data set is useful to researchers studying the relation between habits and practices of expectant mothers and the birth of their children. We will work with a random sample of observations from this data set.

Exploratory analysis

Load the nc data set into our workspace.

```
download.file("http://www.openintro.org/stat/data/nc.RData", destfile = "nc.RData")
load("nc.RData")
```

We have observations on 13 different variables, some categorical and some numerical. The meaning of each variable is as follows.

```
father's age in years.
           fage
                  mother's age in years.
           mage
                  maturity status of mother.
         mature
                  length of pregnancy in weeks.
          weeks
                  whether the birth was classified as premature (premie) or full-term.
         premie
         visits
                  number of hospital visits during pregnancy.
        marital
                  whether mother is married or not married at birth.
                  weight gained by mother during pregnancy in pounds.
         gained
                  weight of the baby at birth in pounds.
         weight
lowbirthweight
                  whether baby was classified as low birthweight (low) or not (not low).
         gender
                  gender of the baby, female or male.
          habit
                  status of the mother as a nonsmoker or a smoker.
       whitemom
                  whether mom is white or not white.
```

Exercise 1 What are the cases in this data set? How many cases are there in our sample?

As a first step in the analysis, we should consider summaries of the data. This can be done using the summary command:

```
summary(nc)
```

As you review the variable summaries, consider which variables are categorical and which are numerical. For numerical variables, are there outliers? If you aren't sure or want to take a closer look at the data, make a graph.

Consider the possible relationship between a mother's smoking habit and the weight of her baby. Plotting the data is a useful first step because it helps us quickly visualize trends, identify strong associations, and develop research questions.

This is a product of OpenIntro that is released under a Creative Commons Attribution-ShareAlike 3.0 Unported (http://creativecommons.org/licenses/by-sa/3.0/). This lab was adapted for OpenIntro by Mine Çetinkaya-Rundel from a lab written by the faculty and TAs of UCLA Statistics.

Exercise 2 Make a side-by-side boxplot of habit and weight. What does the plot highlight about the relationship between these two variables?

The box plots show how the medians of the two distributions compare, but we can also compare the means of the distributions using the following function to split the weight variable into the habit groups, then take the mean of each using the mean function.

```
by(nc$weight, nc$habit, mean)
```

There is an observed difference, but is this difference statistically significant? In order to answer this question we will conduct a hypothesis test.

Inference

Exercise 3 Check if the conditions necessary for inference are satisfied. Note that you will need to obtain sample sizes to check the conditions. You can compute the group size using the same by command above but replacing mean with length.

Exercise 4 Write the hypotheses for testing if the average weights of babies born to smoking and non-smoking mothers are different.

Next, we introduce a new function, inference, that we will use for conducting hypothesis tests and constructing confidence intervals.

Let's pause for a moment to go through the arguments of this custom function.

- The first argument is y, which is the response variable that we are interested in: nc\$weight.
- The second argument is the explanatory variable, x, which is the variable that splits the data into two groups, smokers and non-smokers: nc\$habit.
- The third argument, est, is the parameter we're interested in: "mean" (other options are "median", or "proportion".)
- Next we decide on the type of inference we want: a hypothesis test (ht) or a confidence interval ("ci").
- When performing a hypothesis test, we also need to supply the null value, which in this case is 0, since the null hypothesis sets the two population means equal to each other.
- The alternative hypothesis can be "less", "greater", or "twosided".
- Lastly, the method of inference can be "theoretical" or "simulation" based.

Exercise 5 Change the type argument to "ci" to construct and record a confidence interval for the difference between the weights of babies born to smoking and non-smoking mothers.

By default the function reports an interval for $(\mu_{nonsmoker} - \mu_{smoker})$. We can easily change this order by using the order argument:

On your own

- 1. Calculate a 95% confidence interval for the average length of pregnancies (weeks) and interpret it in context. Note that since you're doing inference on a single population parameter, there is no explanatory variable, so you can omit the x variable from the function.
- 2. Calculate a new confidence interval for the same parameter at the 90% confidence level. You can change the confidence level by adding a new argument to the function: conflevel =0.90.
- 3. Conduct a hypothesis test evaluating whether the average weight gained by younger mothers is different than the average weight gained by mature mothers.
- 4. Now, a non-inference task: Determine the age cutoff for younger and mature mothers. Use a method of your choice, and explain how your method works.
- 5. Pick a pair of numerical and categorical variables and come up with a research question evaluating the relationship between these variables. Formulate the question in a way that it can be answered using a hypothesis test and/or a confidence interval. Answer your question using the inference function, report the statistical results, and also provide an explanation in plain language.
- 6. What concepts from the textbook are covered in this lab? What concepts, if any, are not covered in the textbook? Have you seen these concepts elsewhere, e.g. lecture, discussion section, previous labs, or homework problems? Be specific in your answer.