

# 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.

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

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

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

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

```
nrow(nc)
```

```
## [1] 1000
```

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

```
summary(nc)
```

```
##      fage      mage      mature      weeks
## Min.   :14.00  Min.   :13   mature mom :133  Min.    :20.00
## 1st Qu.:25.00  1st Qu.:22   younger mom:867  1st Qu.:37.00
## Median :30.00  Median :27                      Median :39.00
```

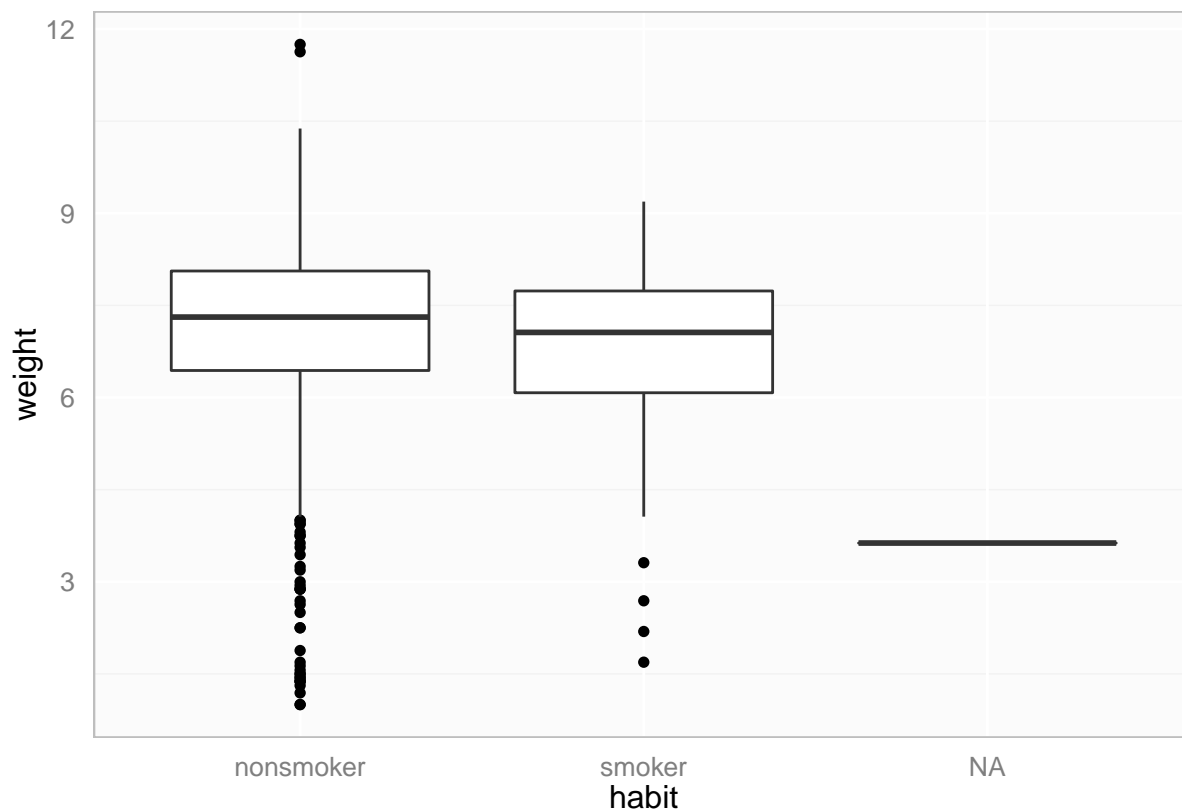
```
## Mean :30.26 Mean :27 Mean :38.33
## 3rd Qu.:35.00 3rd Qu.:32 3rd Qu.:40.00
## Max. :55.00 Max. :50 Max. :45.00
## NA's :171 NA's :2
##      premie      visits      marital      gained
## full term:846 Min. : 0.0 married :386 Min. : 0.00
## premie :152 1st Qu.:10.0 not married:613 1st Qu.:20.00
## NA's : 2 Median :12.0 NA's : 1 Median :30.00
##      Mean :12.1 Mean :30.33
##      3rd Qu.:15.0 3rd Qu.:38.00
##      Max. :30.0 Max. :85.00
##      NA's :9 NA's :27
##      weight lowbirthweight gender habit
## Min. : 1.000 low :111 female:503 nonsmoker:873
## 1st Qu.: 6.380 not low:889 male :497 smoker :126
## Median : 7.310 NA's : 1
## Mean : 7.101
## 3rd Qu.: 8.060
## Max. :11.750
##
##      whitemom
## not white:284
## white :714
## NA's : 2
##
##
##
```

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.

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

```
g1 <- ggplot(data=nc) + geom_boxplot(aes(x=habit, y=weight)) + myTheme
g1
```



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)
```

```
## nc$habit: nonsmoker
## [1] 7.144273
## -----
## nc$habit: smoker
## [1] 6.82873
```

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

## Inference

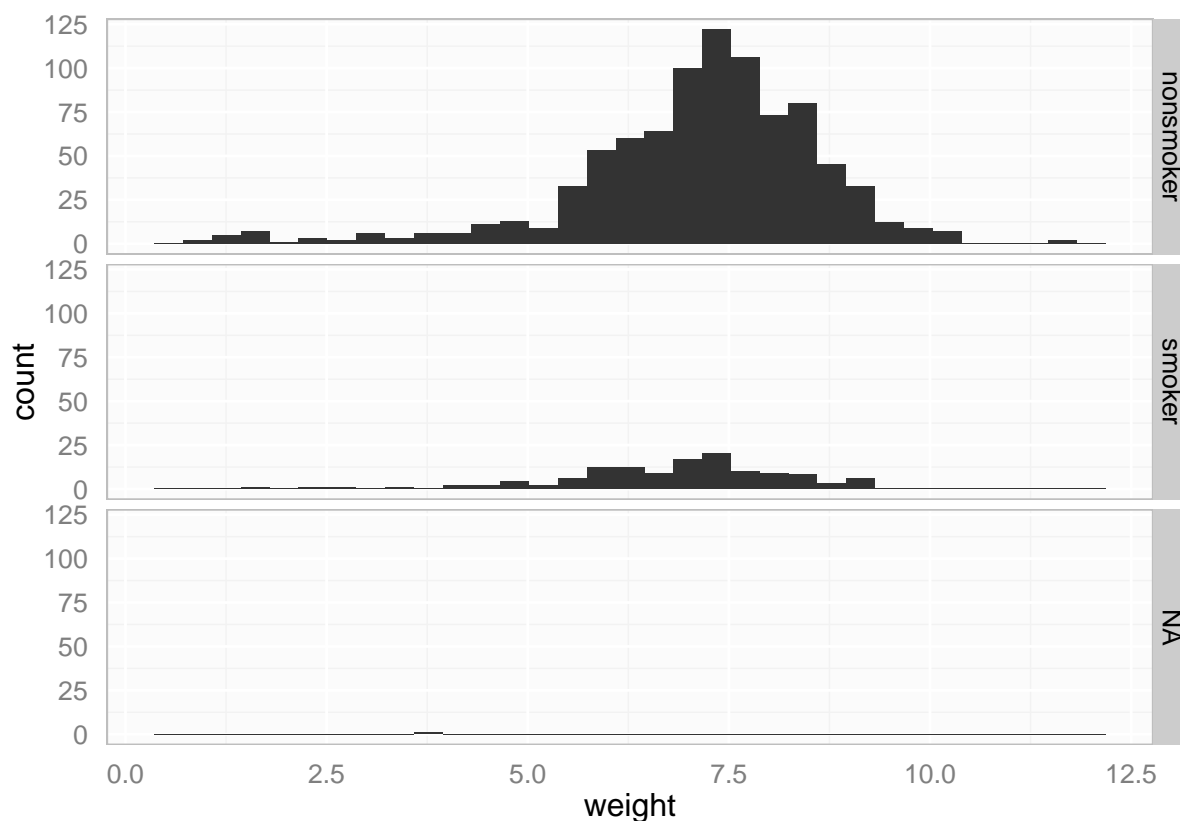
3. Check if the conditions necessary for inference are satisfied. Note that you will need to obtain sample sizes to check the conditions. You can compute the group size using the same `by` command above but replacing `mean` with `length`.

```
# First, how many cases are in each group?
by(nc$weight, nc$habit, length)
```

```
## nc$habit: nonsmoker
## [1] 873
## -----
## nc$habit: smoker
## [1] 126

# Next, show a distribution of weight by habit type.
g1 <- ggplot(data=nc) + geom_bar(aes(x=weight)) + facet_grid(habit ~ .) + myTheme
g1
```

```
## stat_bin: binwidth defaulted to range/30. Use 'binwidth = x' to adjust this.
## stat_bin: binwidth defaulted to range/30. Use 'binwidth = x' to adjust this.
## stat_bin: binwidth defaulted to range/30. Use 'binwidth = x' to adjust this.
```



*The distributions are skewed left, but mostly normal and due to the sample size, this level of skew seems acceptable.*

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

$$H_0 : \mu_s = \mu_{ns}$$

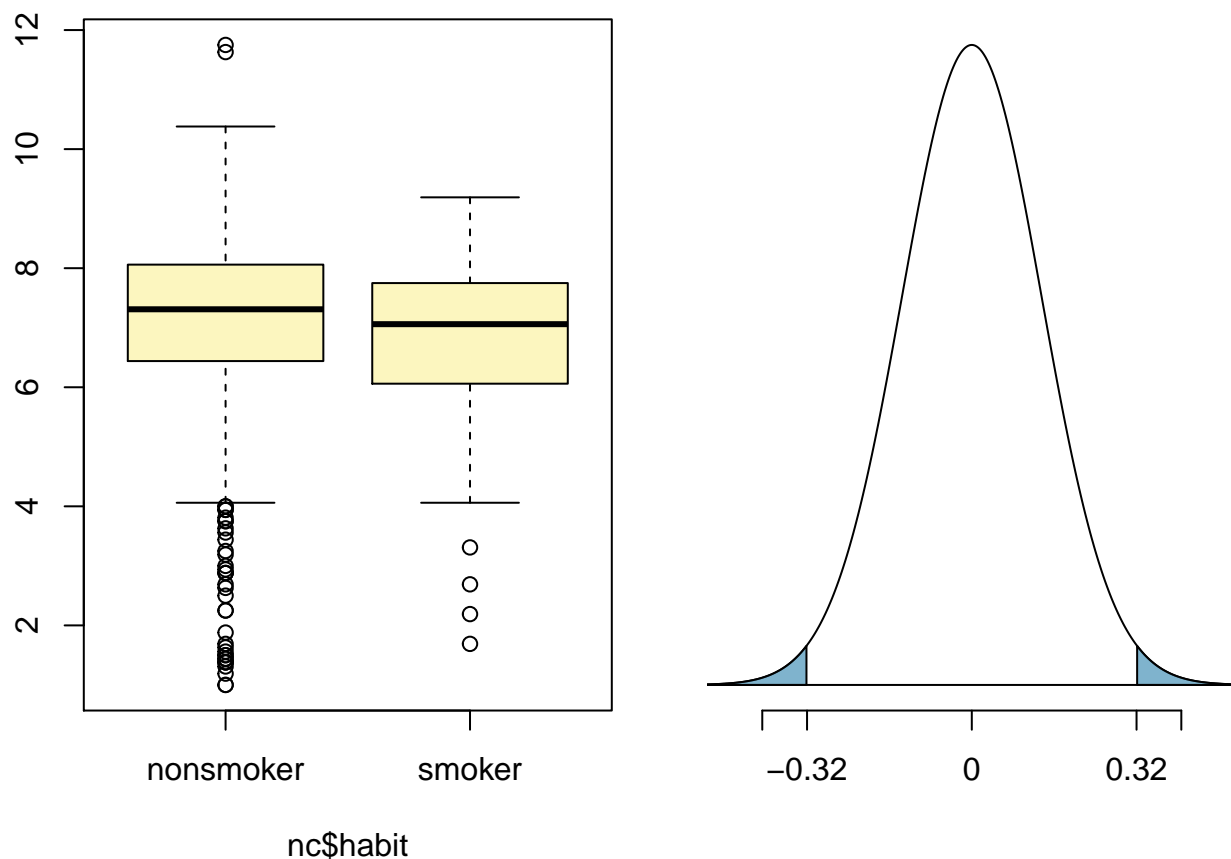
$$H_a : \mu_s \neq \mu_{ns}$$

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

```
inference(y = nc$weight, x = nc$habit, est = "mean", type = "ht", null = 0,
          alternative = "twosided", method = "theoretical")
```

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_nonsmoker = 873, mean_nonsmoker = 7.1443, sd_nonsmoker = 1.5187
## n_smoker = 126, mean_smoker = 6.8287, sd_smoker = 1.3862

## Observed difference between means (nonsmoker-smoker) = 0.3155
##
## H0: mu_nonsmoker - mu_smoker = 0
## HA: mu_nonsmoker - mu_smoker != 0
## Standard error = 0.134
## Test statistic: Z = 2.359
## p-value = 0.0184
```

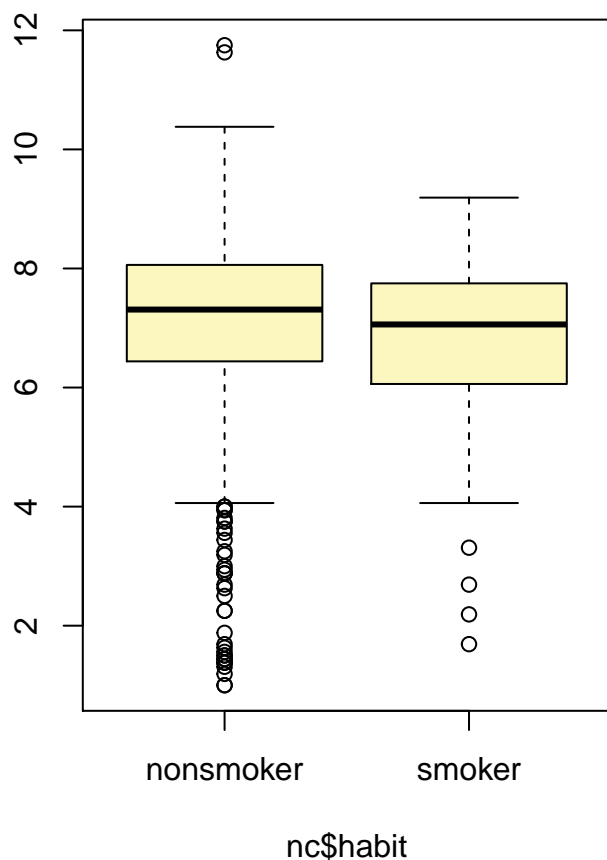


Let's pause for a moment to go through the arguments of this custom function. The first argument is `y`, which is the response variable that we are interested in: `nc$weight`. The second argument is the explanatory variable, `x`, which is the variable that splits the data into two groups, smokers and non-smokers: `nc$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.

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.

```
inference(y = nc$weight, x = nc$habit, est = "mean", type = "ci", null = 0,
          alternative = "twosided", method = "theoretical")
```

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_nonsmoker = 873, mean_nonsmoker = 7.1443, sd_nonsmoker = 1.5187
## n_smoker = 126, mean_smoker = 6.8287, sd_smoker = 1.3862
```

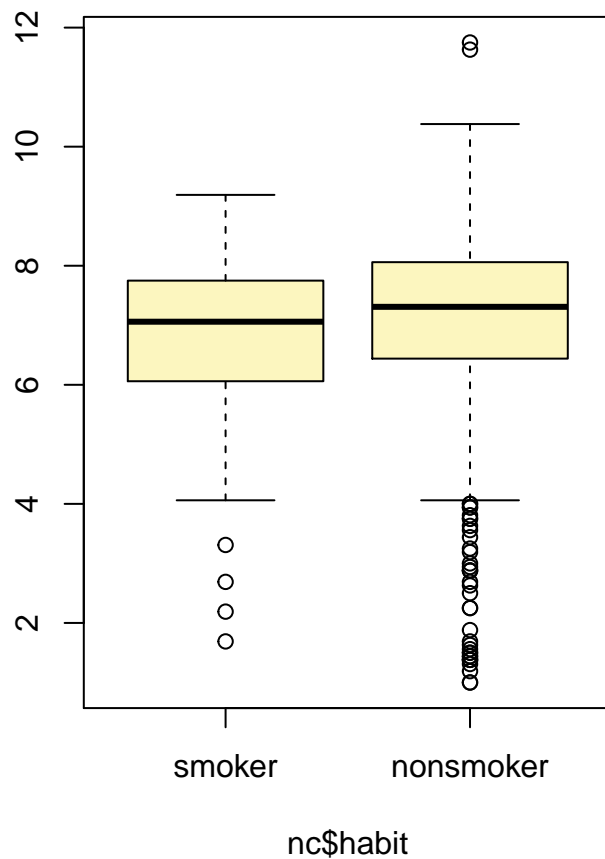


```
## Observed difference between means (nonsmoker-smoker) = 0.3155
##
## Standard error = 0.1338
## 95 % Confidence interval = ( 0.0534 , 0.5777 )
```

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

```
inference(y = nc$weight, x = nc$habit, est = "mean", type = "ci", null = 0,
          alternative = "twosided", method = "theoretical",
          order = c("smoker", "nonsmoker"))
```

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_smoker = 126, mean_smoker = 6.8287, sd_smoker = 1.3862
## n_nonsmoker = 873, mean_nonsmoker = 7.1443, sd_nonsmoker = 1.5187
```



```
## Observed difference between means (smoker-nonsmoker) = -0.3155
##
## Standard error = 0.1338
## 95 % Confidence interval = ( -0.5777 , -0.0534 )
```

## On your own

- 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.

[Docs for inference function](#)

```
meanWeeks <- mean(nc$weeks, na.rm=TRUE)
meanWeeks
```

```
## [1] 38.33467
```

```
sdWeeks <- sd(nc$weeks, na.rm=TRUE)
sdWeeks
```

```
## [1] 2.931553
```

```
n <- nrow(nc$weeks)
n
```

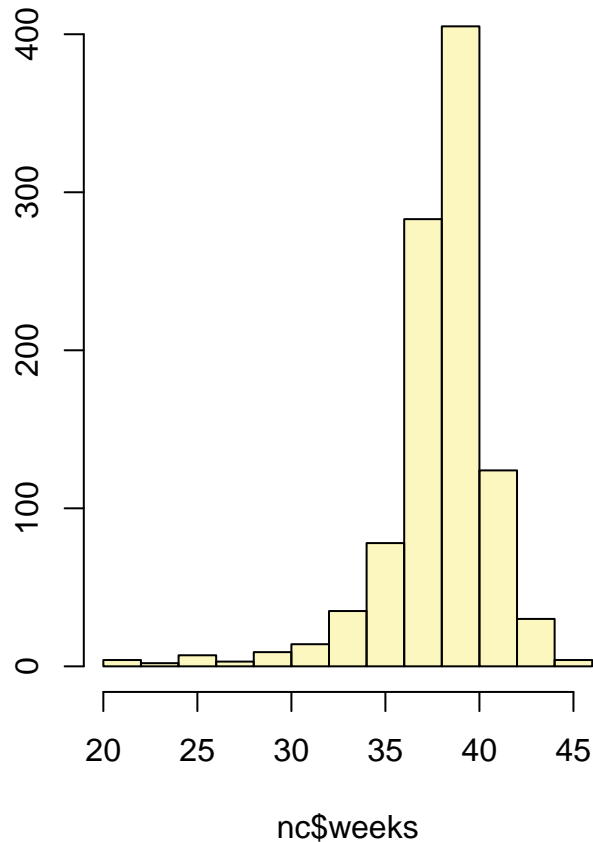
```
## NULL
```

```
seWeeks <- sqrt(sdWeeks / n)
seWeeks
```

```
## numeric(0)
```

```
# Calling the inference function.
inference(y=nc$weeks, est="mean", type="ci", method="theoretical")
```

```
## Single mean
## Summary statistics:
```





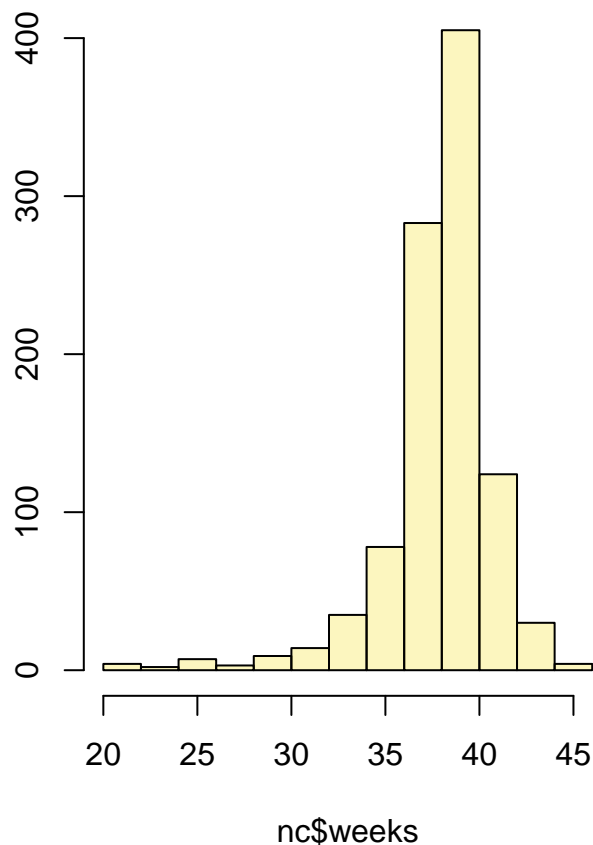
```
## mean = 38.3347 ; sd = 2.9316 ; n = 998
## Standard error = 0.0928
## 95 % Confidence interval = ( 38.1528 , 38.5165 )
```

*I am 95% confident that the population mean weeks is between 38.1528 and 38.516.*

- 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: `conflvel = 0.90`.

```
inference(y=nc$weeks, est="mean", type="ci", method="theoretical", conflvel=0.90)
```

```
## Single mean
## Summary statistics:
```



```
## mean = 38.3347 ; sd = 2.9316 ; n = 998
## Standard error = 0.0928
## 90 % Confidence interval = ( 38.182 , 38.4873 )
```

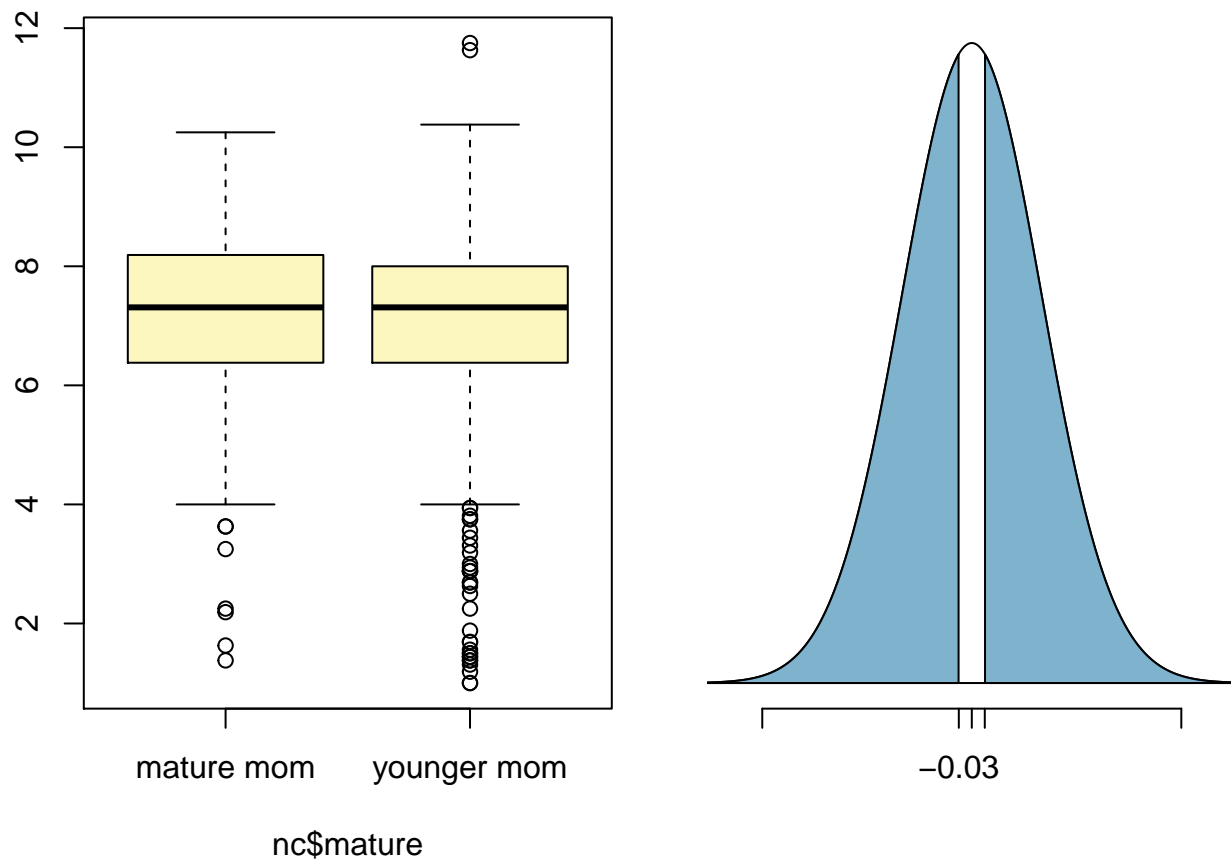
*I am 90% confident that the population mean weeks is between 38.182 and 38.4873.*

- Conduct a hypothesis test evaluating whether the average weight gained by younger mothers is different than the average weight gained by mature mothers.

```
inference(y=nc$weight, x=nc$mature, type="ht", est="mean",
          null=0, method="theoretical", alternative="twosided")
```

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_mature mom = 133, mean_mature mom = 7.1256, sd_mature mom = 1.6591
## n_younger mom = 867, mean_younger mom = 7.0972, sd_younger mom = 1.4855

## Observed difference between means (mature mom-younger mom) = 0.0283
##
## H0: mu_mature mom - mu_younger mom = 0
## HA: mu_mature mom - mu_younger mom != 0
## Standard error = 0.152
## Test statistic: Z = 0.186
## p-value = 0.8526
```



*Due to the p-value being so high, at 0.8526, I fail to reject the null hypothesis. Therefore I conclude the weight of babies born to younger moms is not different than born to mature moms.*

- 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.

*The maximum “younger mom” age would close to the cutoff, and the minimum “mature mom” age would be close to the cutoff. In our case, younger mom max is 34 and mature mom age is 35, so it seems the cutoff for younger is less than or equal to 34, and mature is greater than 34.*

```
max(nc[nc$mature == "younger mom"],$mage)
```

```
## [1] 34
```

```
min(nc[nc$mature == "mature mom"],$mage)
```

```
## [1] 35
```

- 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.

*Does the number of average visits to the doctor differ between married and unmarried moms?*

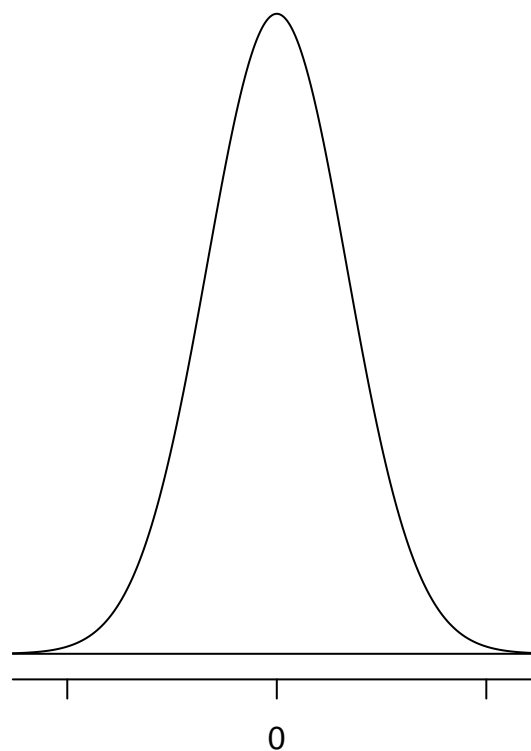
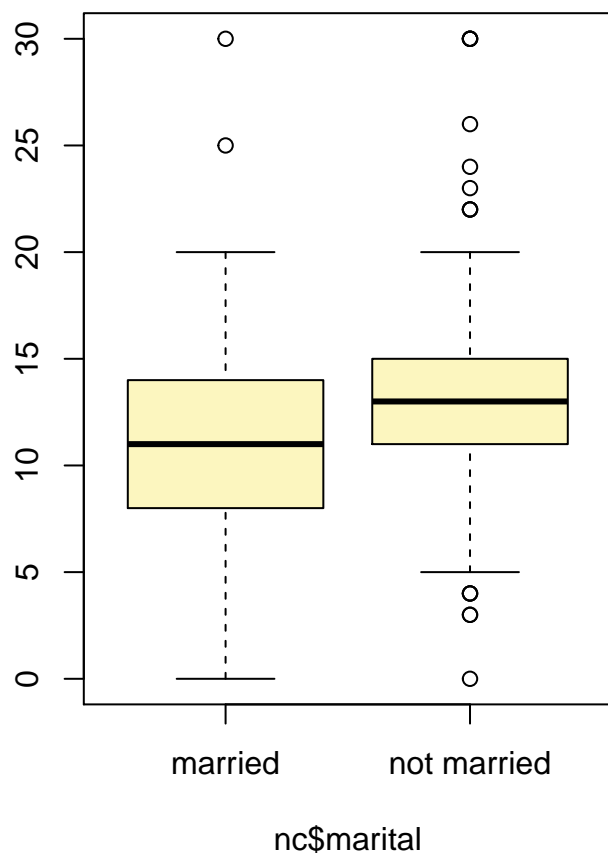
$H_0 : \mu_{vm} = \mu_{vn}$

$H_a : \mu_{vm} \neq \mu_{vn}$

```
inference(y=nc$visits, x=nc$marital, type="ht", est="mean",
          null=0, method="theoretical", alternative="twosided")
```

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_married = 380, mean_married = 10.9553, sd_married = 4.2408
## n_not married = 611, mean_not married = 12.82, sd_not married = 3.5883

## Observed difference between means (married-not married) = -1.8647
##
## H0: mu_married - mu_not married = 0
## HA: mu_married - mu_not married != 0
## Standard error = 0.262
## Test statistic: Z = -7.13
## p-value = 0
```



*Given the  $p\text{-value}=0$ , I reject the null hypothesis and conclude that the average number of visits to the doctor is different between married and non-married moms.*

This is a product of OpenIntro that is released under a [Creative Commons Attribution-ShareAlike 3.0 Unported](#). This lab was adapted for OpenIntro by Mine Çetinkaya-Rundel from a lab written by the faculty and TAs of UCLA Statistics.