

# 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.

variable	description
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.

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

Each case in this dataset are birth details for babies born in North Carolina. There are 1,000 cases in this 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)
```

```
##      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
```

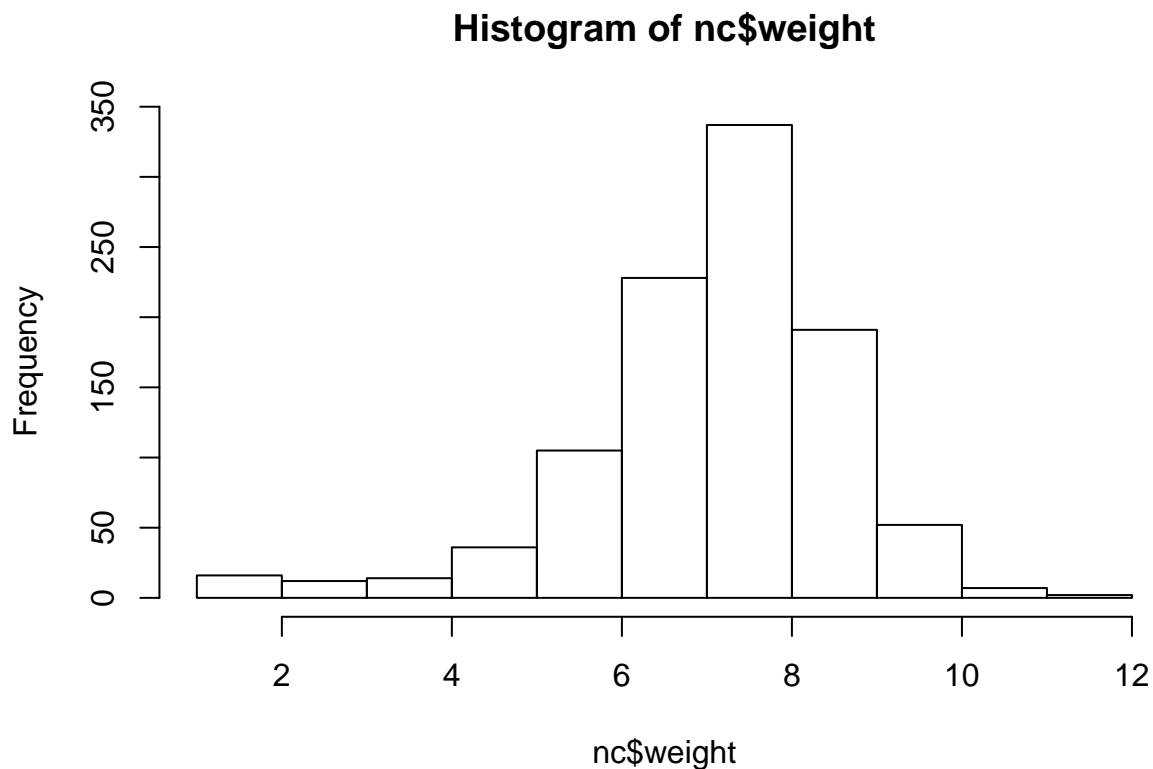
```
## not white:284
## white :714
## NA's : 2
##
##
##
##
```

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.

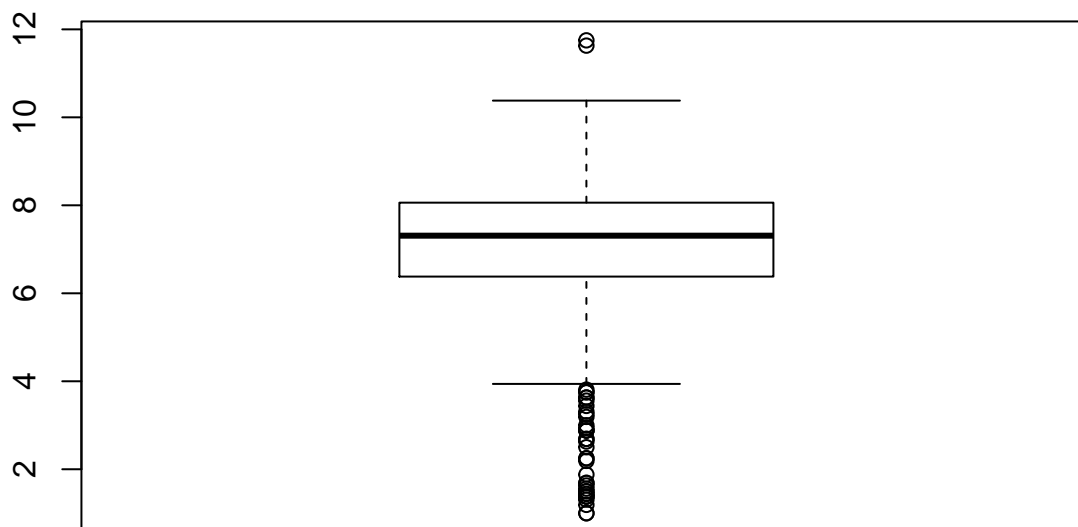
It looks like 'mature', 'premie', 'marital', 'lowbirthweight', 'gender', 'habit, and 'whitemom' are categorical variables and 'fage', 'mage', 'weeks', 'visits', 'gained', and 'weight' are numerical variables.

However, it would be helpful to check to see if there are any outliers for 'weight', since it is such an important factor on the baby's health:

```
hist(nc$weight)
```



```
boxplot(nc$weight)
```



```
summary(nc$weight)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1.000   6.380   7.310   7.101   8.060  11.750
```

We can see that there are a fair amount of outliers for a baby's weight – those that have a registered birth weight of under about 4 pounds is considered uncommon, and a few babies have registered birth weights above 10.5 pounds, which is also considered uncommon in this sample. The boxplot shows these outliers.

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?

Here's the side-by-side boxplot of 'habit' and 'weight':

```
boxplot(weight~habit,data=nc, main="Relation Between Mother's Smoking Habits and Baby's Weight",
        ylab="Baby's Weight", xlab="Mother Smoker/Non-Smoker")
```

## Relation Between Mother's Smoking Habits and Baby's Weight



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`.

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

```
## nc$habit: nonsmoker
```

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

It appears that the samples are independent from one another given that this is a random sample (outlined earlier).

Additionally, it appears that there is normality of this sample, given that  $n >$  or equal to 30 for both non-smoking and smoking mothers in this sample, as indicated from the test above.

Therefore, both conditions for inference are satisfied.

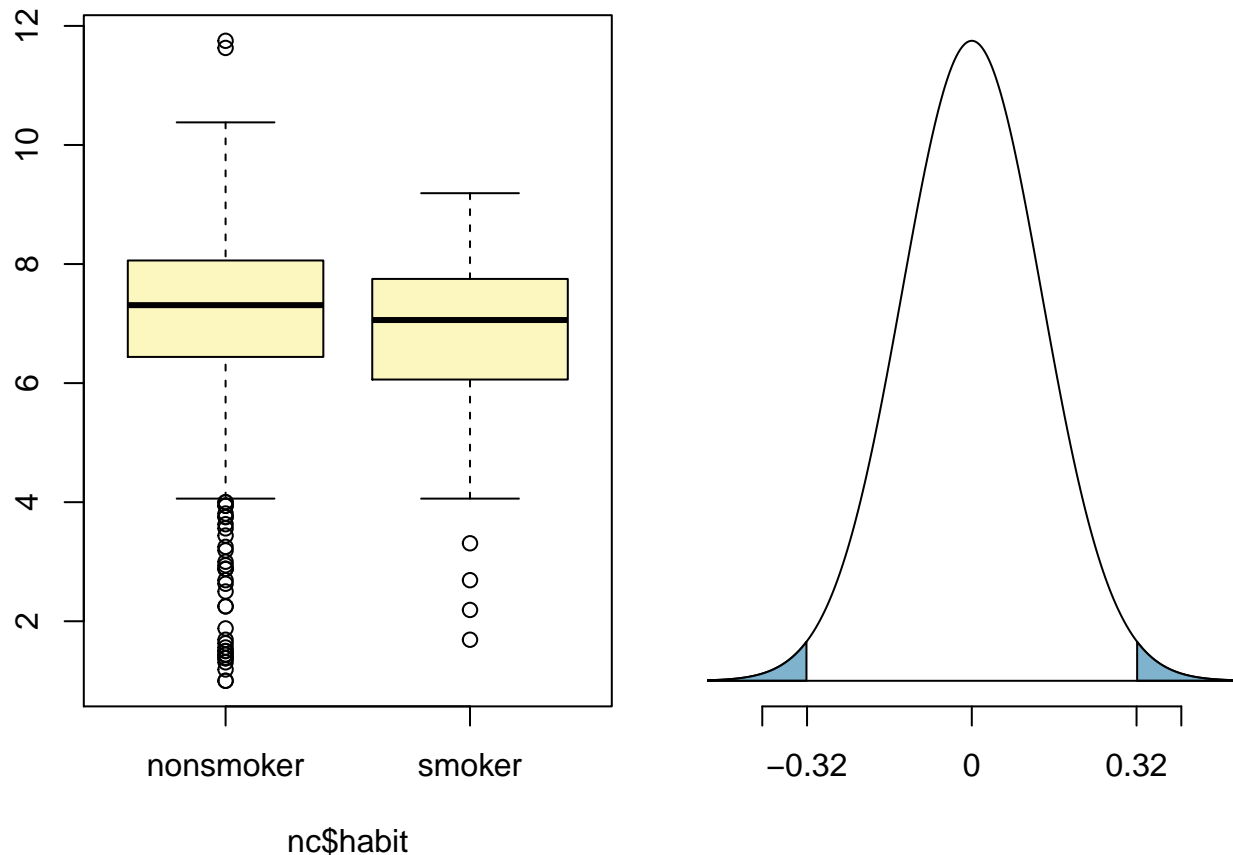
4. Write the hypotheses for testing if the average weights of babies born to smoking and non-smoking mothers are different.
- **H0:** There is no difference between the average weights of babies born to smoking and non-smoking mothers.
  - **HA:** There is a difference between the average weights of babies born to smoking and non-smoking mothers.

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
```

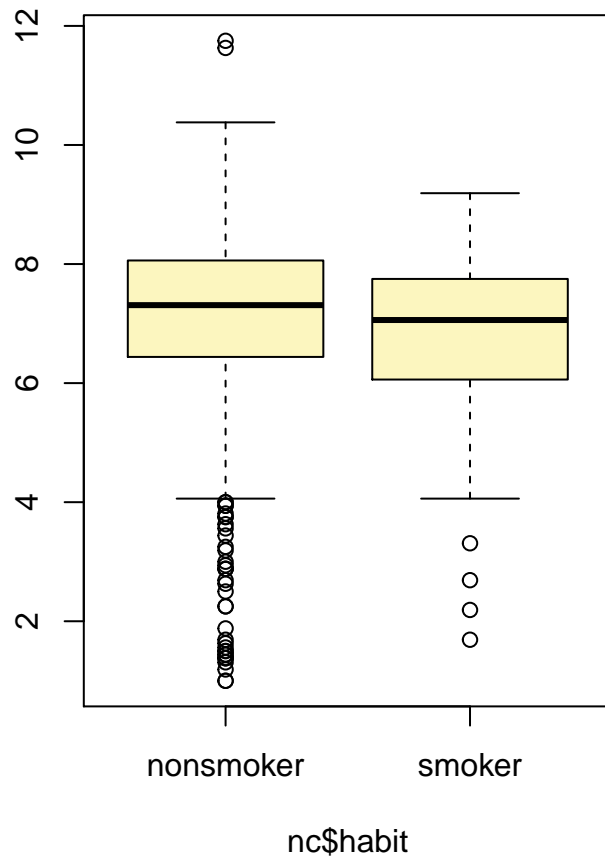


Let's pause for a moment to go through the arguments of this custom function. The first argument is `y`, which is the response variable that we are interested in: `nc$weight`. The second argument is the explanatory variable, `x`, which is the variable that splits the data into two groups, smokers and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

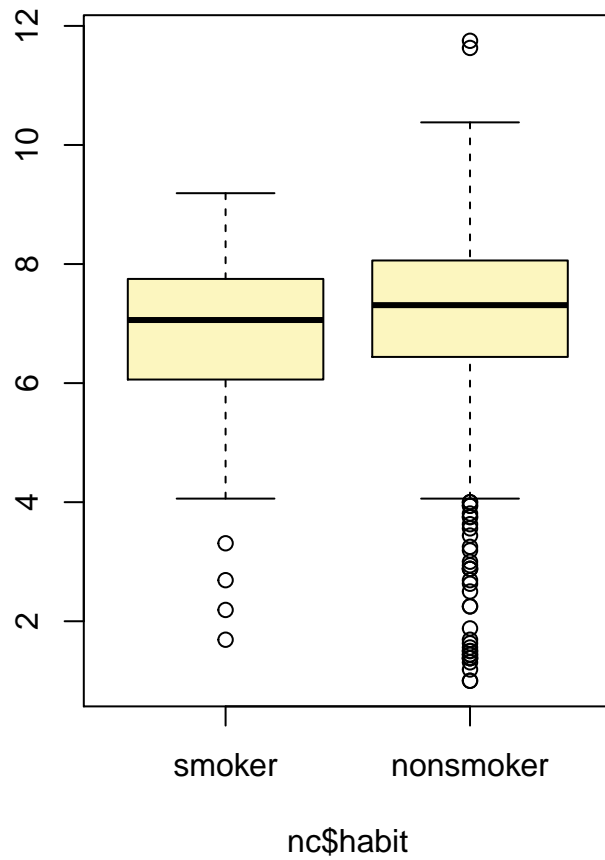
The confidence interval is (0.0534, 0.5777).

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 )
```

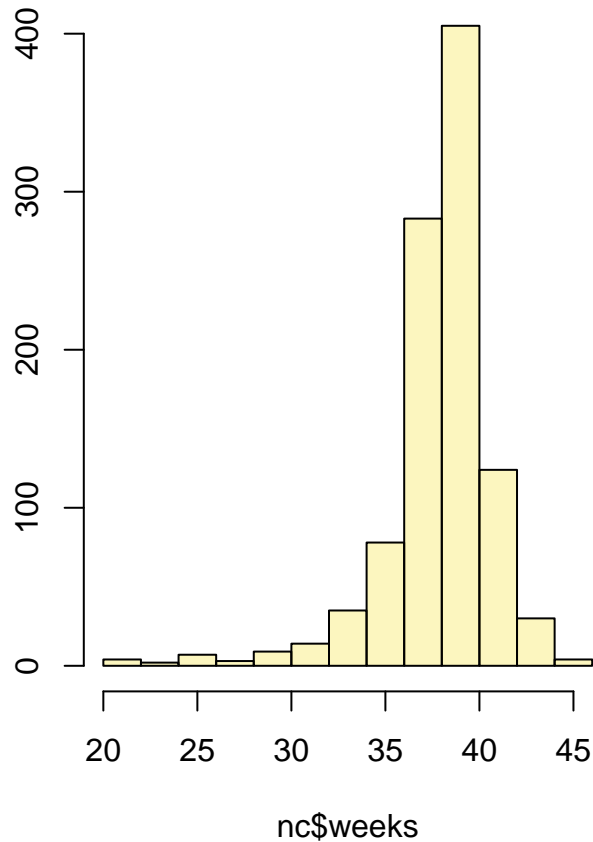
---

### 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.

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

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



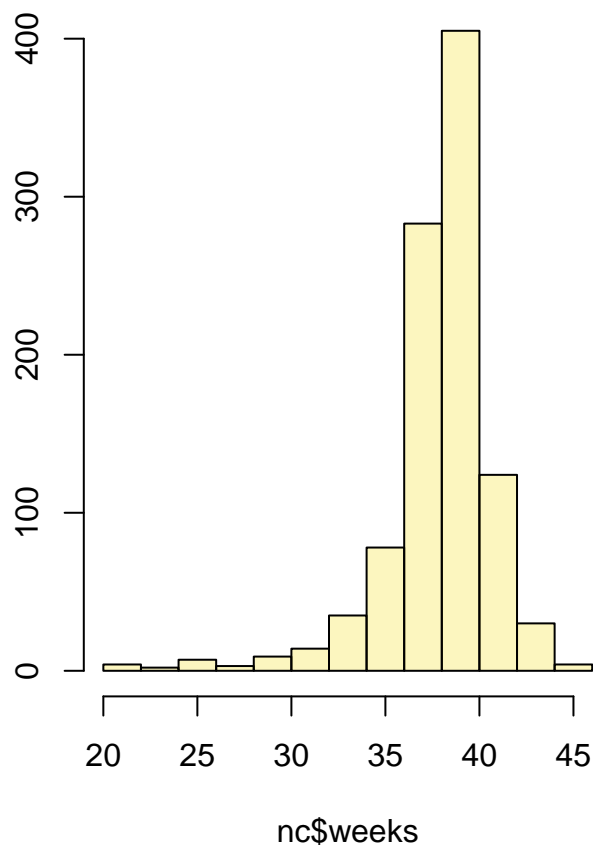
```
## mean = 38.3347 ; sd = 2.9316 ; n = 998
## Standard error = 0.0928
## 95 % Confidence interval = ( 38.1528 , 38.5165 )
```

The 95% confidence interval for the average length of pregnancies for this sample (in weeks) is: **38.15 weeks to 38.52 weeks.**

- 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`.

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

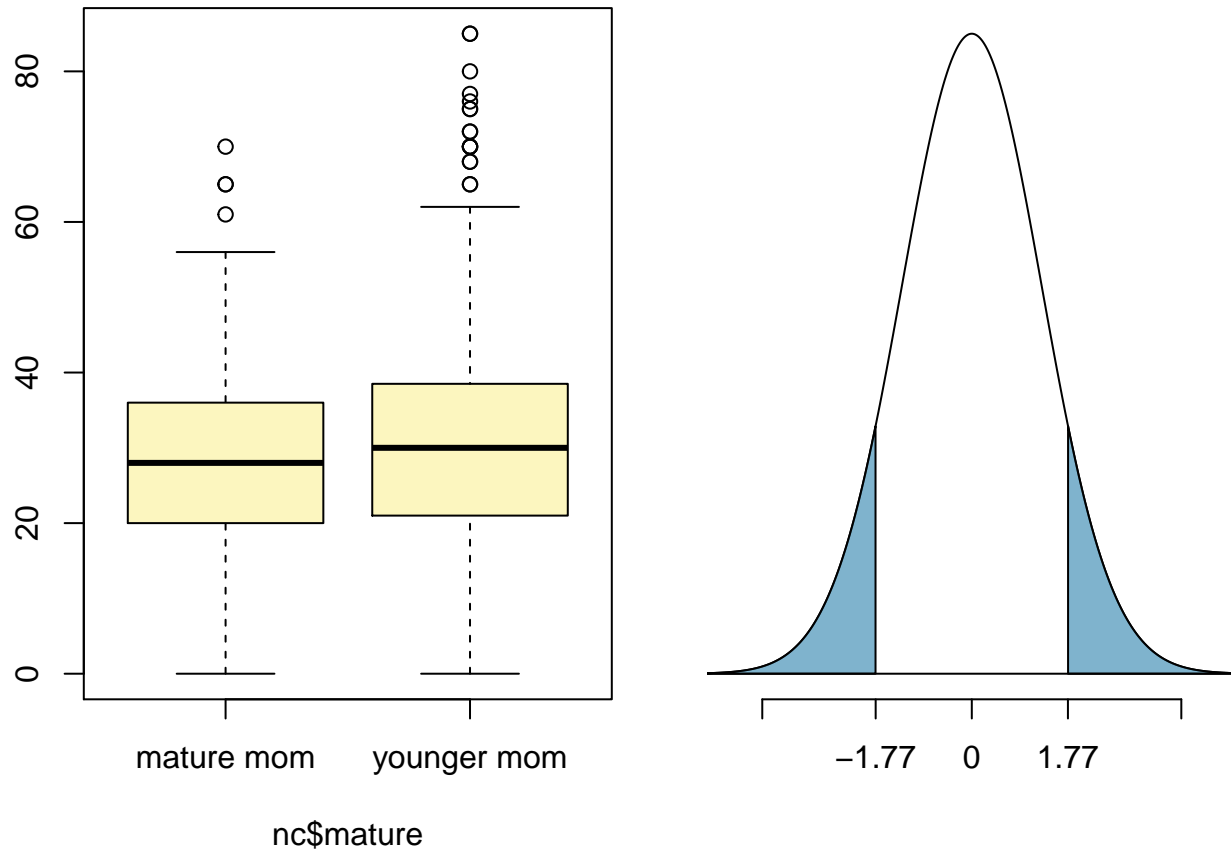
The 90% confidence interval for the average length of pregnancies for this sample (in weeks) is: 38.18 weeks to 38.49 weeks.

- Conduct a hypothesis test evaluating whether the average weight gained by younger mothers is different than the average weight gained by mature mothers.
- **H0:** There is no difference between the average weight gained by younger mothers and the average weight gained by mature mothers.
- **HA:** There is a difference between the average weight gained by younger mothers and the average weight gained by mature mothers.

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

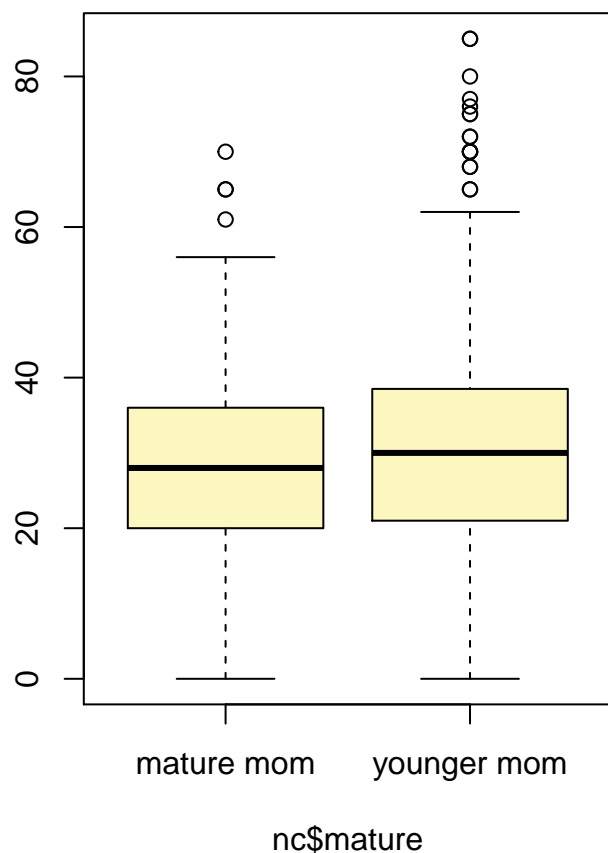
```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_mature mom = 129, mean_mature mom = 28.7907, sd_mature mom = 13.4824
## n_younger mom = 844, mean_younger mom = 30.5604, sd_younger mom = 14.3469
```

```
## Observed difference between means (mature mom-younger mom) = -1.7697
##
## H0: mu_mature mom - mu_younger mom = 0
## HA: mu_mature mom - mu_younger mom != 0
## Standard error = 1.286
## Test statistic: Z = -1.376
## p-value = 0.1686
```



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

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_mature mom = 129, mean_mature mom = 28.7907, sd_mature mom = 13.4824
## n_younger mom = 844, mean_younger mom = 30.5604, sd_younger mom = 14.3469
```



```
## Observed difference between means (mature mom-younger mom) = -1.7697
##
## Standard error = 1.2857
## 95 % Confidence interval = ( -4.2896 , 0.7502 )
```

According to the hypothesis test, it appears that there is no difference between the average weight gained by younger mothers and the average weight gained by mature mothers. With a p-value = 0.1686, and a confidence interval that spans 0, we can accept the null hypothesis.

- 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.

```
library(tidyverse)
```

```
## -- Attaching packages -----
## v ggplot2 3.2.1    v purrr  0.3.2
## v tibble  2.1.3    v dplyr  0.8.3
## v tidyr   0.8.3    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.4.0

## -- Conflicts -----
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
nc %>%
  group_by(mature, mage) %>%
  summarise(count = n()) %>%
  arrange(mage) %>%
  summarise(maxage = max(mage), minage = min(mage)) %>%
  gather(x, n, 2:3) %>%
  select(mature, n) %>%
  filter((max(n) & mature == 'younger mom') | (min(n) & mature == 'mature mom')) %>%
  arrange(mature, -n) %>%
  filter(!row_number() == 1 & !row_number() == n())
```

```
## # A tibble: 2 x 2
##   mature      n
##   <fct>    <int>
## 1 mature mom    35
## 2 younger mom   34
```

Using tidyverse and some data wrangling methods, I've found that the cutoff is between 34 and 35 years of age for those considered to be a younger mother from mature mothers.

I was able to use `group_by` and summary methods in order to find the min and max age broken down by the categorical variable of `mature`, then finding the min and max values of those groupings, I was able to isolate the min age for mature moms and the max age for younger moms.

- 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.

**Research question:** Is there a difference in the average weight of a baby depending on whether or not a mother is considered a younger mother or a mature mother?

To answer this question, we'll assume the following hypotheses:

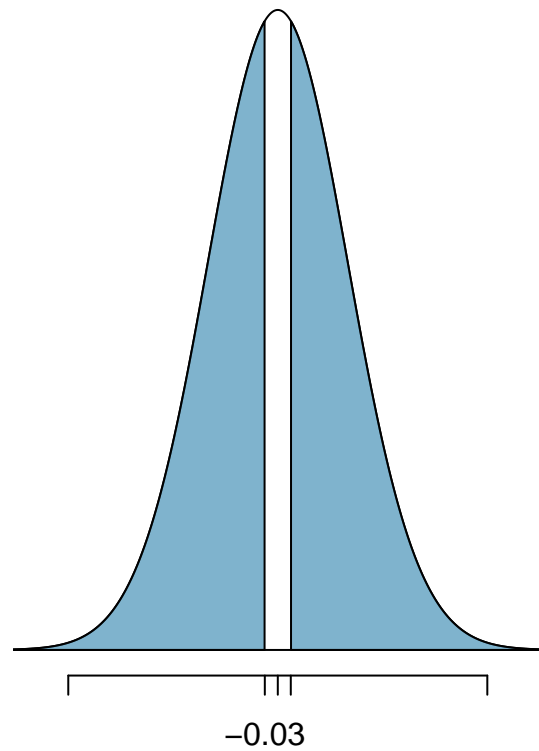
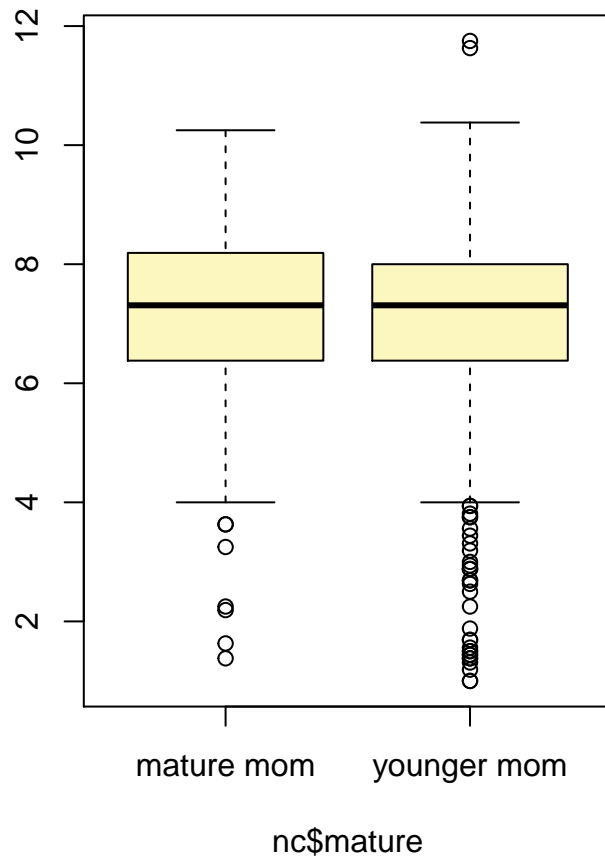
- **H0:** There is no difference between the average weights of babies born from younger mothers and the average weights of a babies born from mature mothers.
- **HA:** There is a difference between the average weights of babies born from younger mothers and the average weights of a babies born from mature mothers.

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

```
## 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

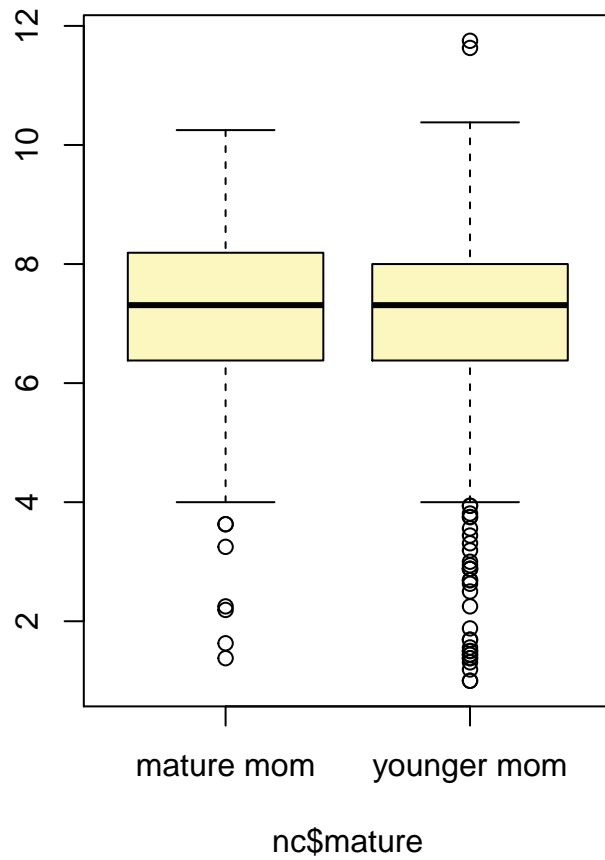
## Observed difference between means (mature mom-younger mom) = 0.0283
##
```

```
## 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
```



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

```
## 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
##
## Standard error = 0.1525
## 95 % Confidence interval = ( -0.2705 , 0.3271 )
```

We accept the null hypothesis that there is no difference between the average weights of babies born from younger mothers and the average weights of babies born from mature mothers. Since the p-value = 0.186, and the 95% confidence interval spans 0, we can accept the null hypothesis.

In plain language, it appears that the age of the mother bearing a child will not have an impact on their baby's birth weight.