

Inference for numerical data

Subhalaxmi Rout

22/03/2020

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 (<code>premie</code>) 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?

Answer

A case is a single birth in the state of North Carolina. There are total 1000 cases in this dataset.

```
dim(nc)
```

```
## [1] 1000 13
```

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      premie
## Min.   :14.00  Min.   :13    mature mom :133  Min.   :20.00  full term:846
## 1st Qu.:25.00  1st Qu.:22    younger mom:867  1st Qu.:37.00  premie   :152
## Median :30.00  Median :27                                Median :39.00  NA's     : 2
## 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
##      visits      marital      gained      weight
## Min.   : 0.0    married   :386  Min.   : 0.00  Min.   : 1.000
## 1st Qu.:10.0    not married:613  1st Qu.:20.00  1st Qu.: 6.380
## Median :12.0    NA's       : 1    Median :30.00  Median : 7.310
## Mean   :12.1                                Mean   :30.33  Mean   : 7.101
## 3rd Qu.:15.0                                3rd Qu.:38.00  3rd Qu.: 8.060
## Max.   :30.0                                Max.   :85.00  Max.   :11.750
## NA's   :9                                NA's   :27
## lowbirthweight  gender      habit      whitemom
## low           :111  female:503  nonsmoker:873  not white:284
## not low:889    male  :497  smoker  :126  white   :714
##                                     NA's     : 1  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.

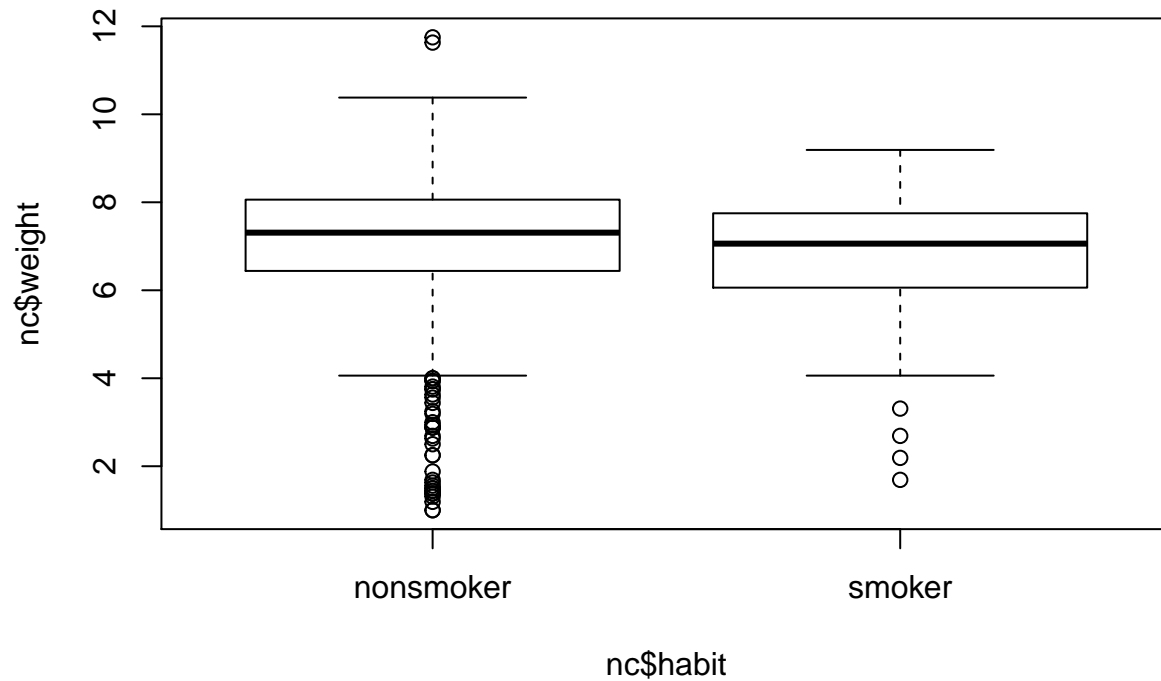
Consider the possible relationship between a mother's smoking habit and the weight of her baby. Plotting the data is a useful first step because it helps us quickly visualize trends, identify strong associations, and develop research questions.

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

Answer

The boxplot shows that the median birth weight of newborns of mothers who is non-smoker is higher than the newborns mothers who is smoker.

```
boxplot(nc$weight ~ nc$habit)
```



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`.

Answer

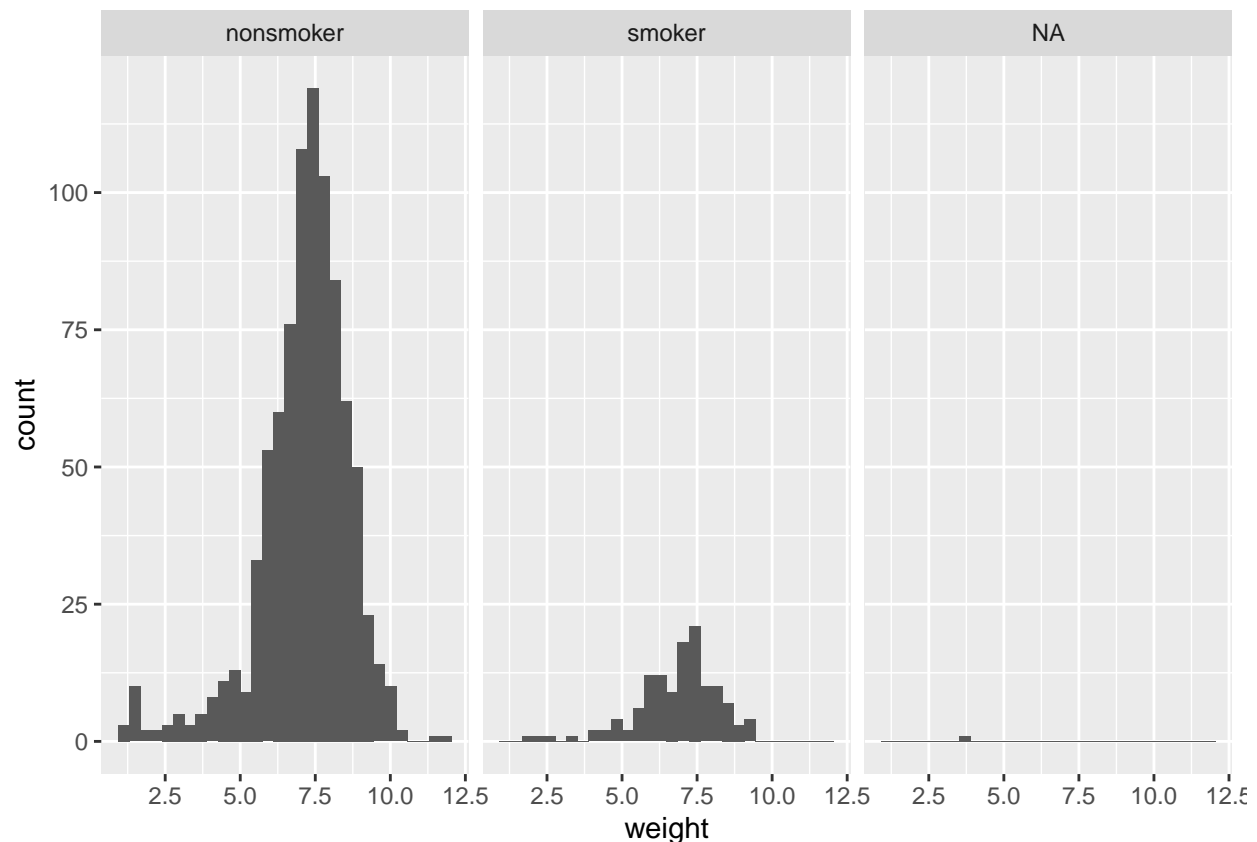
```
by(nc$weight, nc$habit, length)
```

```
## nc$habit: nonsmoker
## [1] 873
```

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

```
library(ggplot2)
ggplot(nc,aes(x=weight)) + geom_histogram() + facet_grid(~nc$habit)
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



samples are approximately normally distributed in both group. The sample size in each group is > 30 . The sample is definitely less than 10% birth in North Carolina. Observations in each group seems independent. So from this we can say this is a normal distribution.

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

Answer

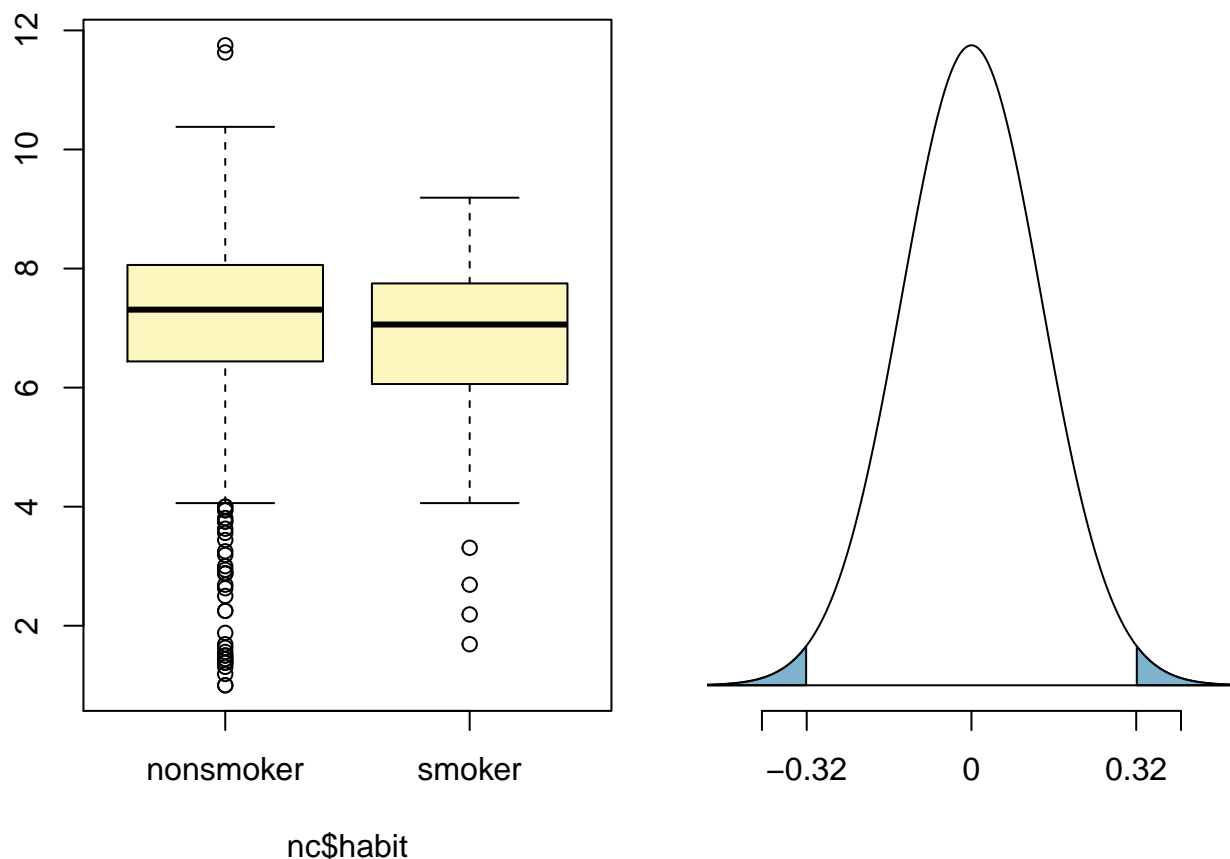
H_0 : average weights of babies born to smoking and non-smoking mother are same H_A : average weights of babies born to smoking and non-smoking mother are not same

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

```
DATA606::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
```



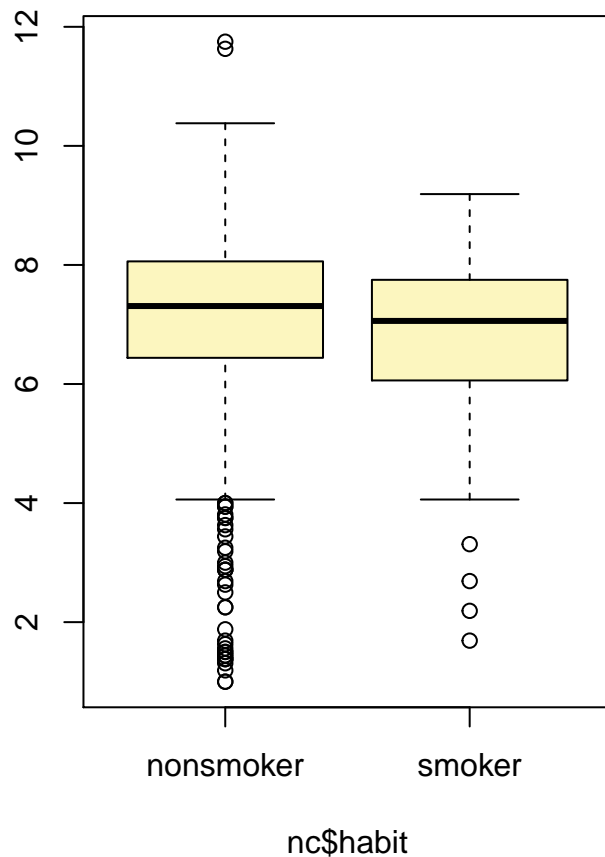
Let's pause for a moment to go through the arguments of this custom function. The first argument is `y`, which is the response variable that we are interested in: `nc$weight`. The second argument is the explanatory variable, `x`, which is the variable that splits the data into two groups, smokers and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

Answer

```
DATA606::inference( nc$weight, nc$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
```



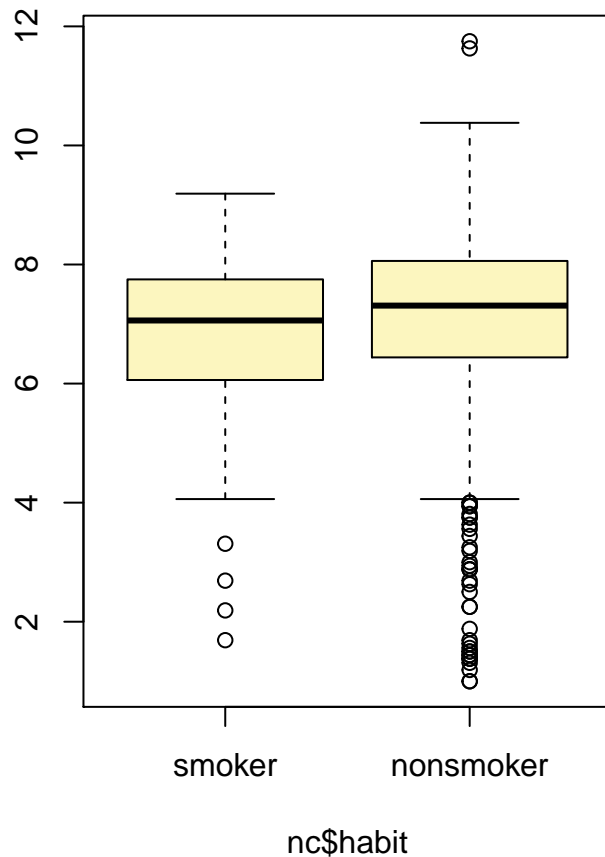
```
## Observed difference between means (nonsmoker-smoker) = 0.3155  
##  
## Standard error = 0.1338  
## 95 % Confidence interval = ( 0.0534 , 0.5777 )
```

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

```
DATA606::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 )
```

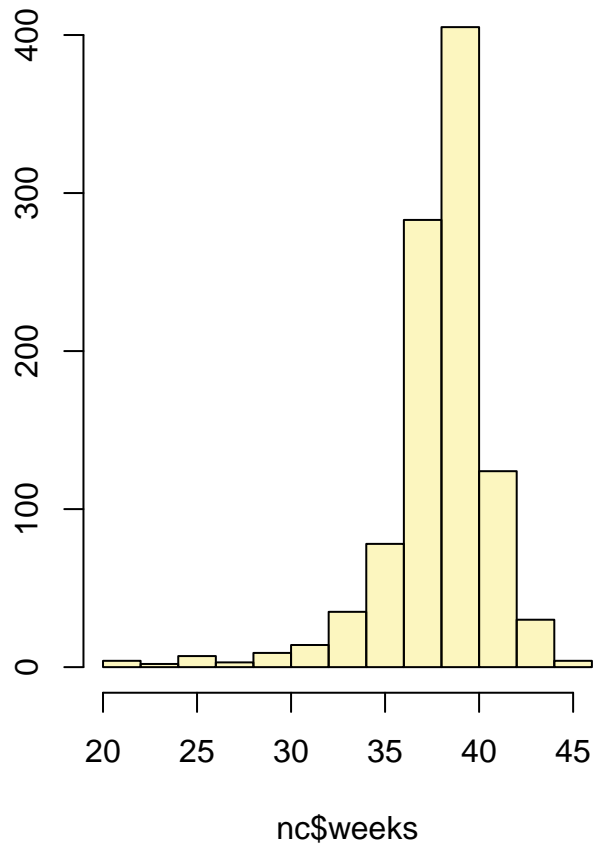
On your own

- 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.

Answer

```
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 )
```

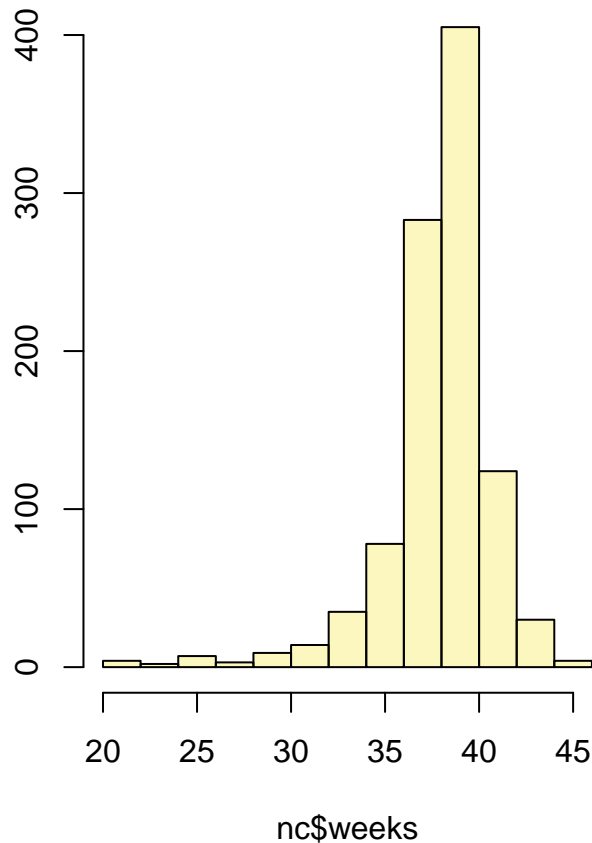
95 % Confidence interval is (38.1528 , 38.5165)

- 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`.

Answer

```
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 )
```

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

Answer

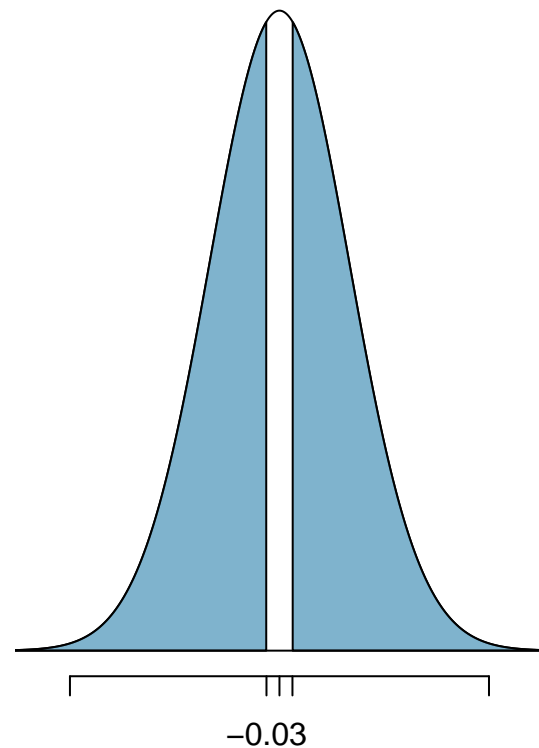
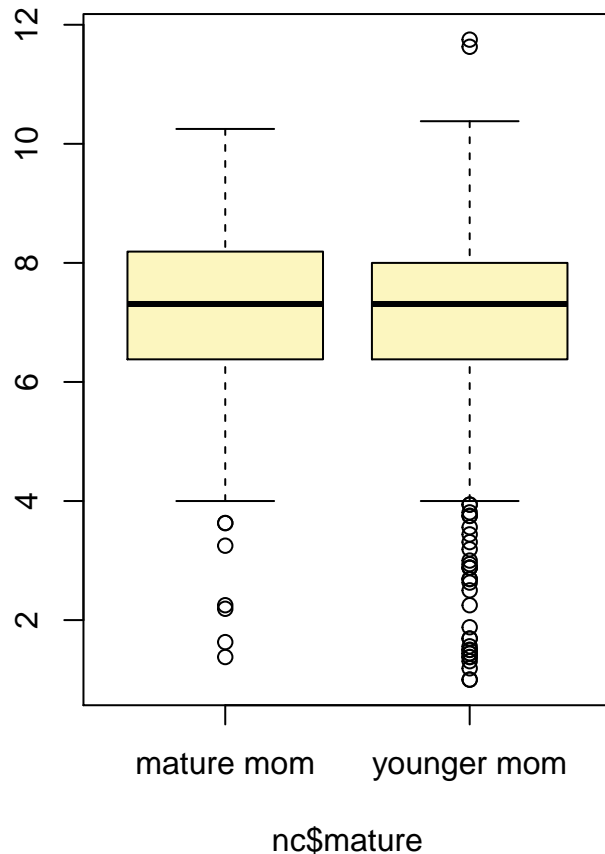
H0: There is no difference between average weight gained by mature mothers and younger mothers. HA: There is difference between average weight gained by mature mothes and younger mothers.

```
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
```



Since the confidence interval $(-4.2896, 0.7502)$. Based on above data, we accept reject the null hypothesis and we are saying that there is no difference in birth weight of babies born to younger and mature mothers.

- 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.

Answer

```
by(nc$age, nc$mature, range)
```

```
## nc$mature: mature mom
## [1] 35 50
## -----
## nc$mature: younger mom
## [1] 13 34
```

From the above analysis, we can see the age for younger mother is between 13 to 34. The age for mature mother is between 35 to 50. Here I have used `by()` to calculate age consideration for mature and younger mother.

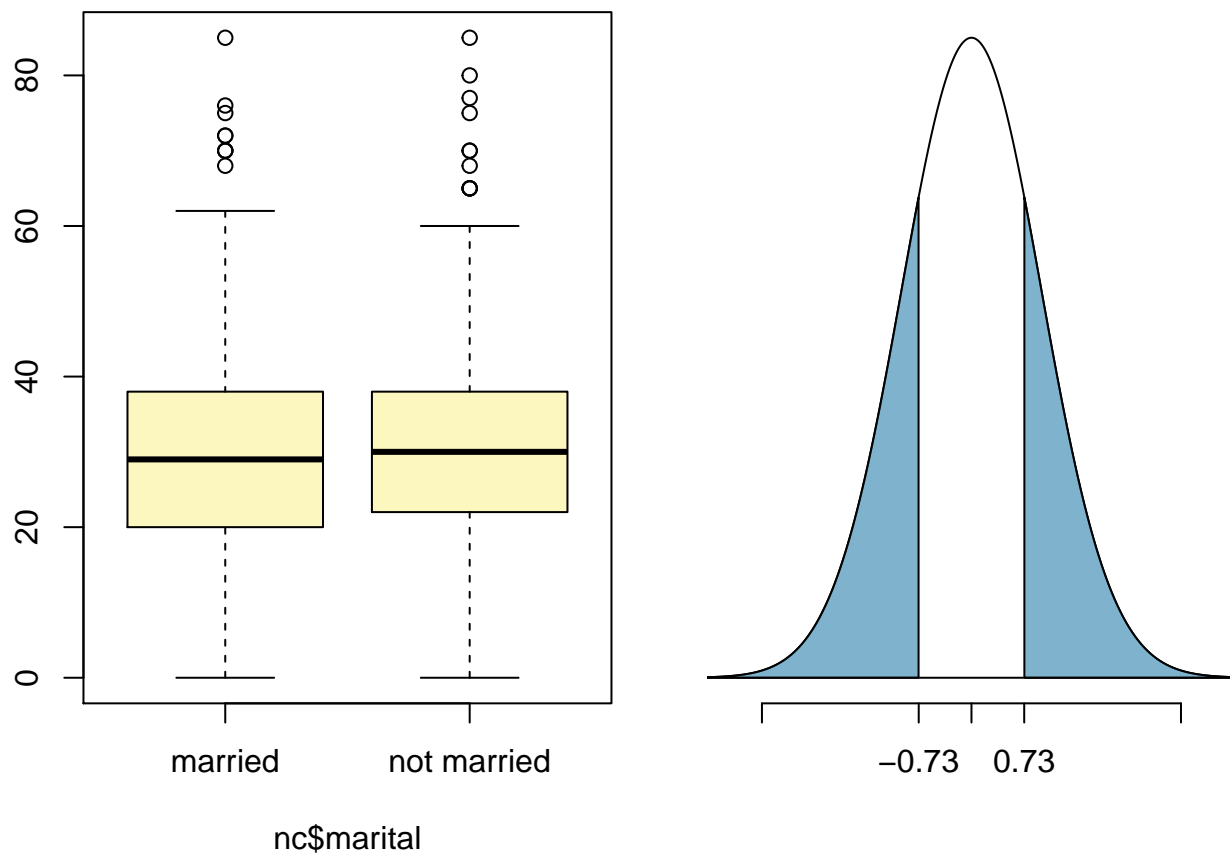
- Pick a pair of numerical and categorical variables and come up with a research question evaluating the relationship between these variables. Formulate the question in a way that it can be answered using a hypothesis test and/or a confidence interval. Answer your question using the `inference` function, report the statistical results, and also provide an explanation in plain language.

Answer

Hypothesis test: H0: There is no difference in the mean of the weight gained during pregnancy of between married and unmarried mothers. HA: There is difference in the mean of the weight gained during pregnancy of between married and unmarried mothers

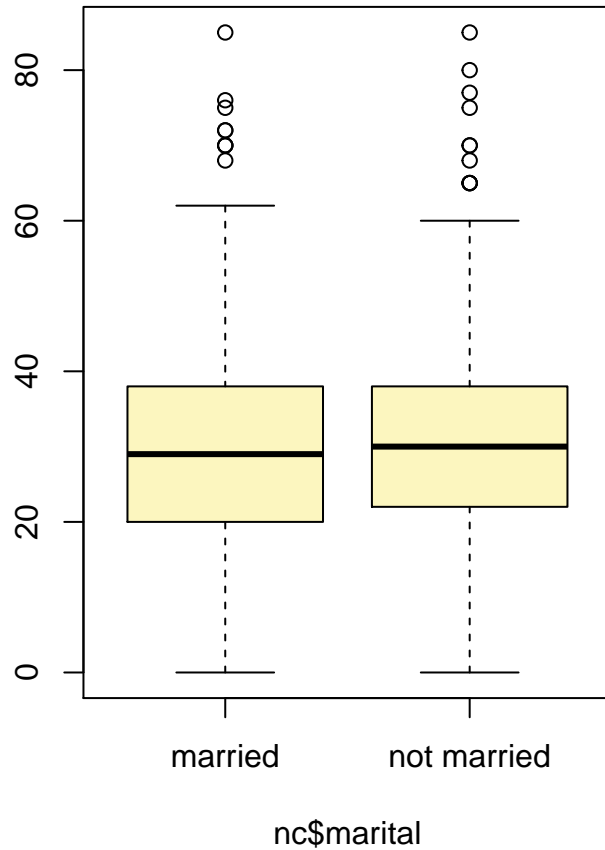
```
inference(y = nc$gained, x = nc$marital, est = "mean", type = "ht", null = 0,  
          alternative = "twosided", method = "theoretical")
```

```
## Response variable: numerical, Explanatory variable: categorical  
## Difference between two means  
## Summary statistics:  
## n_married = 370, mean_married = 29.873, sd_married = 15.2721  
## n_not married = 603, mean_not married = 30.6036, sd_not married = 13.5757  
  
## Observed difference between means (married-not married) = -0.7307  
##  
## H0:  $\mu_{\text{married}} - \mu_{\text{not married}} = 0$   
## HA:  $\mu_{\text{married}} - \mu_{\text{not married}} \neq 0$   
## Standard error = 0.967  
## Test statistic:  $Z = -0.755$   
## p-value = 0.4502
```



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

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_married = 370, mean_married = 29.873, sd_married = 15.2721
## n_not married = 603, mean_not married = 30.6036, sd_not married = 13.5757
```



```
## Observed difference between means (married-not married) = -0.7307
##
## Standard error = 0.9675
## 95 % Confidence interval = ( -2.6269 , 1.1655 )
```

From the above data, we cannot reject the hypothesis. There is no evidence based on statistical data to show that there is a difference between the weight gained by married and unmarried mothers during pregnancy.