# 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 |
|---|---|
| fage | father's age in years. |
| mage | mother's age in years. |
| mature | maturity status of mother. |
| weeks | length of pregnancy in weeks. |
| premie | whether the birth was classified as premature (premie) or full-term. |
| visits | number of hospital visits during pregnancy. |
| marital | whether mother is `married` or `not married` at birth. |
| gained | weight gained by mother during pregnancy in pounds. |
| weight | weight of the baby at birth in pounds. |
| 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? JR Answer: There are 1000 cases of births in NC

```
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
##
##
##
##
```

```
tail(nc)
```

```
##       fage mage     mature weeks    premie visits      marital gained weight
## 995    NA   41 mature mom    33    premie     13 not married      0   5.69
## 996    47   42 mature mom    40 full term     10 not married     26   8.44
## 997    34   42 mature mom    38 full term     18 not married     20   6.19
## 998    39   45 mature mom    40 full term     15 not married     32   6.94
## 999    55   46 mature mom    31    premie      8 not married     25   4.56
## 1000   45   50 mature mom    39 full term     14 not married     23   7.13
##      lowbirthweight gender      habit  whitemom
## 995          not low female nonsmoker not white
## 996          not low   male nonsmoker not white
## 997          not low female nonsmoker    white
## 998          not low female nonsmoker    white
## 999              low female nonsmoker not white
## 1000         not low female nonsmoker    white
```

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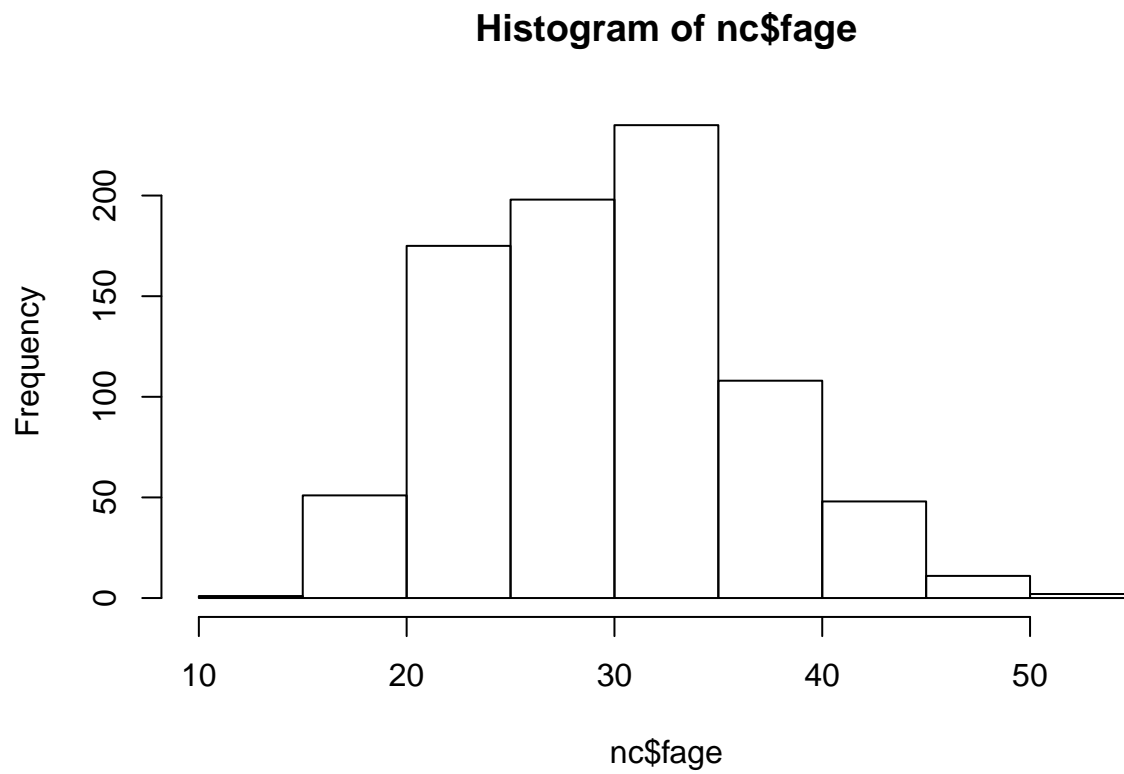
Consider the possible relationship between a mother's smoking habit and the weight of her baby. Plotting the data is a useful first step because it helps us quickly visualize trends, identify strong associations, and develop research questions.

```
str(nc)
```

```
## 'data.frame':    1000 obs. of  13 variables:
##  $ fage          : int  NA NA 19 21 NA NA 18 17 NA 20 ...
##  $ mage          : int  13 14 15 15 15 15 15 15 16 16 ...
##  $ mature        : Factor w/ 2 levels "mature mom","younger mom": 2 2 2 2 2 2 2 2 2 2 ...
##  $ weeks         : int  39 42 37 41 39 38 37 35 38 37 ...
##  $ premie        : Factor w/ 2 levels "full term","premie": 1 1 1 1 1 1 1 2 1 1 ...
##  $ visits        : int  10 15 11 6 9 19 12 5 9 13 ...
##  $ marital       : Factor w/ 2 levels "married","not married": 1 1 1 1 1 1 1 1 1 1 ...
##  $ gained        : int  38 20 38 34 27 22 76 15 NA 52 ...
##  $ weight        : num  7.63 7.88 6.63 8 6.38 5.38 8.44 4.69 8.81 6.94 ...
##  $ lowbirthweight: Factor w/ 2 levels "low","not low": 2 2 2 2 2 1 2 1 2 2 ...
##  $ gender        : Factor w/ 2 levels "female","male": 2 2 1 2 1 2 2 2 2 1 ...
##  $ habit         : Factor w/ 2 levels "nonsmoker","smoker": 1 1 1 1 1 1 1 1 1 1 ...
##  $ whitemom      : Factor w/ 2 levels "not white","white": 1 1 2 2 1 1 1 1 2 2 ...
```
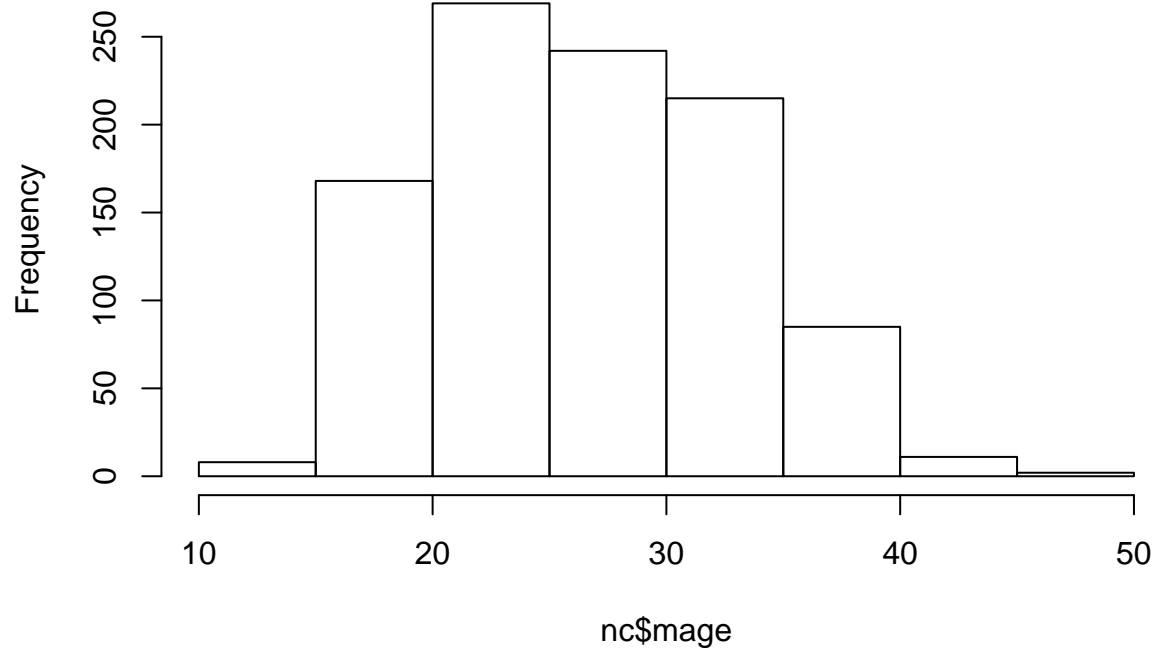
JR Answer: In the weight histogram - outliers are 1 lbs and 11.85 lbs int and num are numberical: fage, mage, weeks,visits, gained, weight factors are categorical hist($fage)
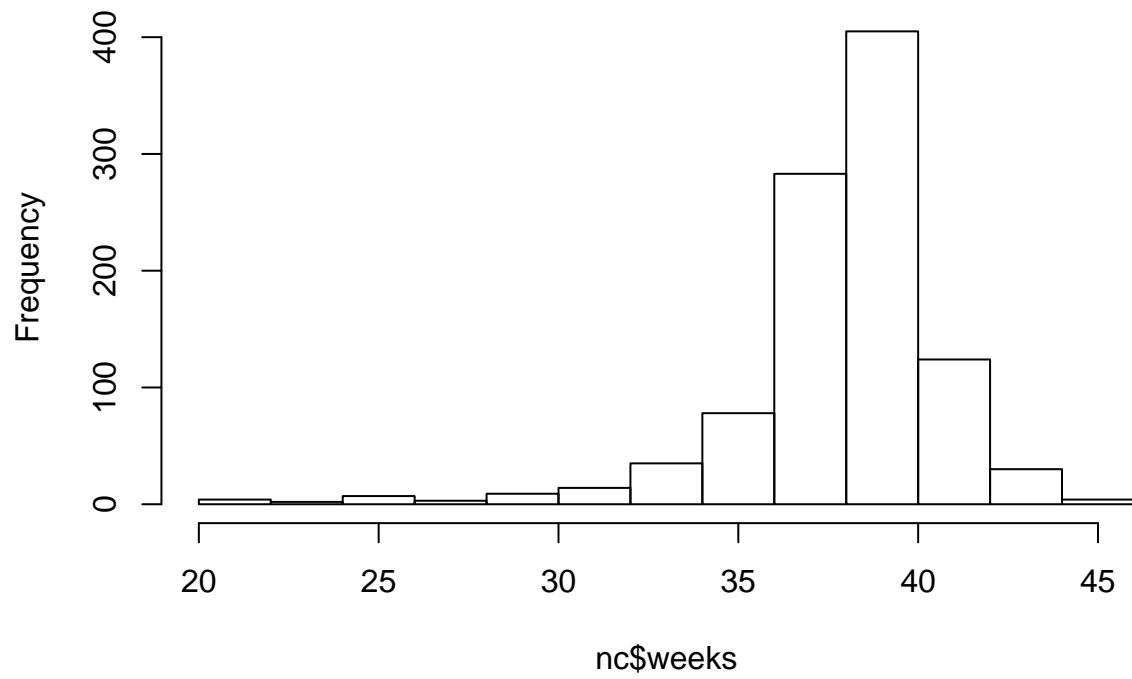
```r
hist(nc$fage)
```
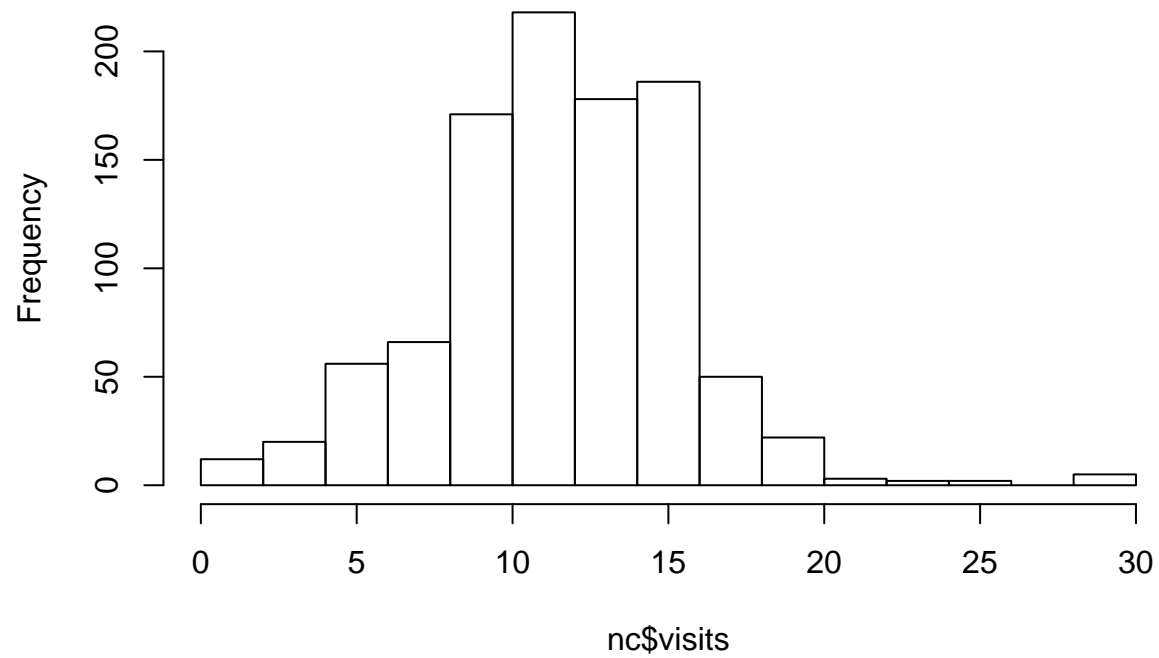
## Histogram of nc$fage



```r
hist(nc$mage)
```

**Histogram of nc$mage**
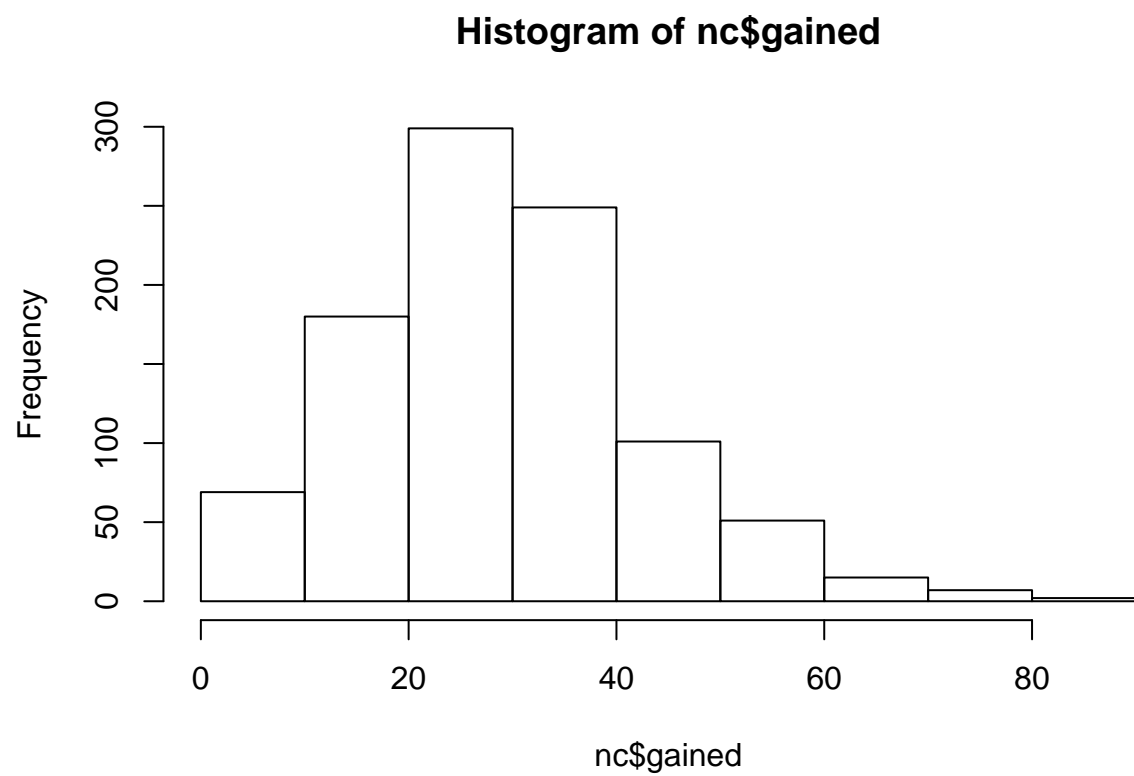


```
hist(nc$weeks)
```

**Histogram of nc$weeks**

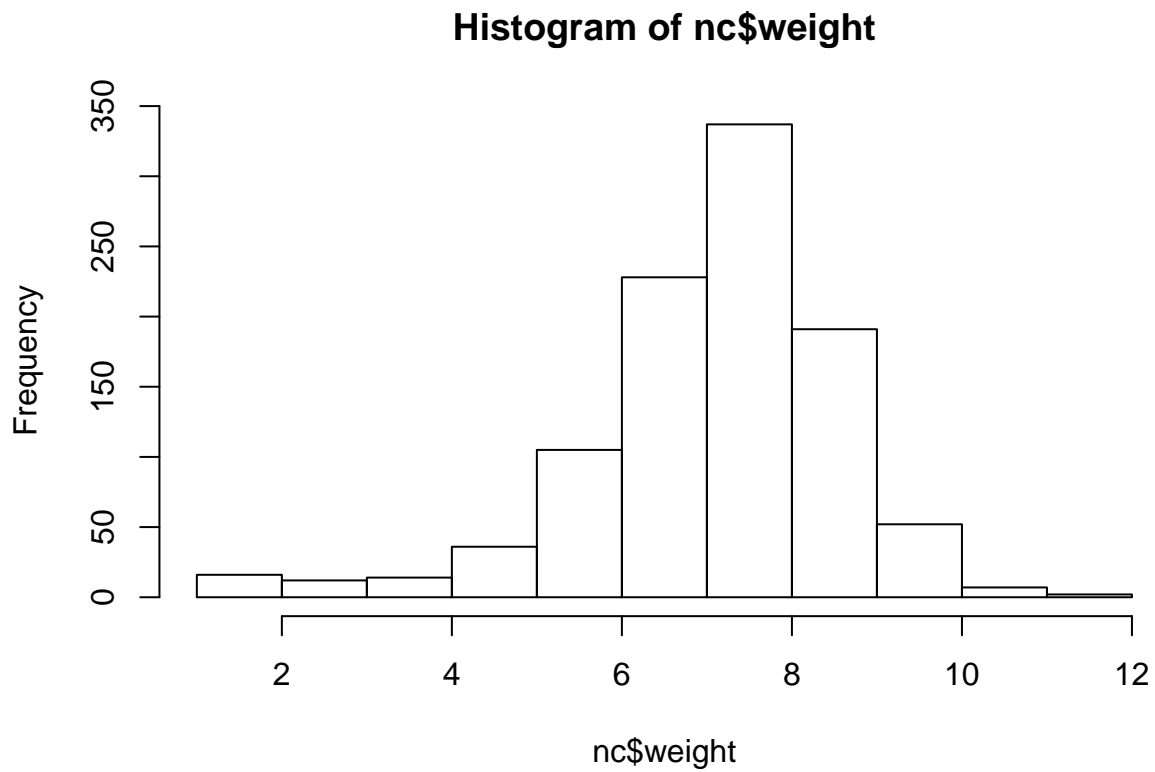Frequency

```
hist(nc$visits)
```

## Histogram of nc$visits



nc$visits

```r
hist(nc$gained)
```

**Histogram of nc$gained**



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

## Histogram of nc$weight



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

```r
# Boxplot of habit and weight
boxplot(weight~habit,data=nc, main="Mother's Habit vs Baby's Weight",
    ylab="Baby Weight", xlab="Mother Smoker/Non-Smoker")
```

## Mother's Habit vs 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
```

JR Answer: Sample observations are independent as are the sameple groups. The sample sizes are less then 10% of the population size. Sample size is sufficient as to not worry about skew.
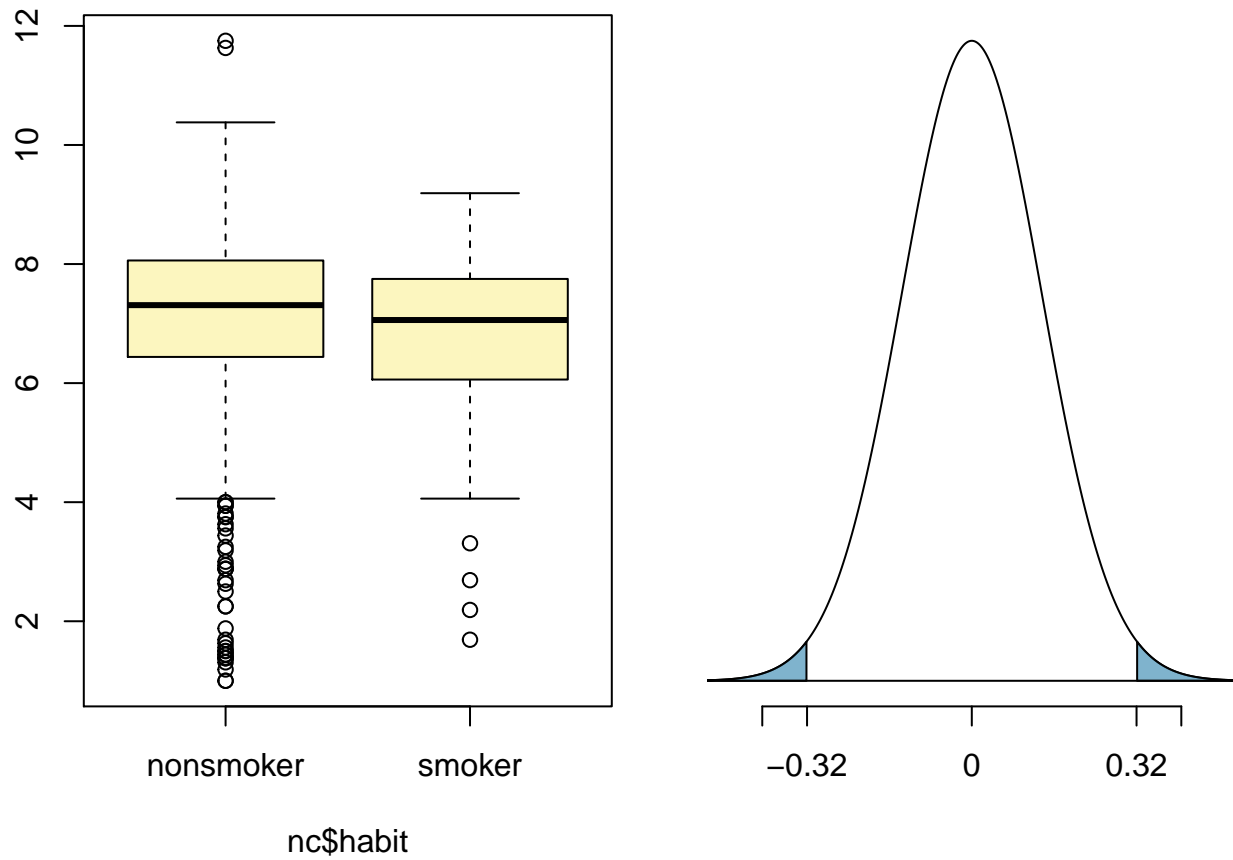
4. Write the hypotheses for testing if the average weights of babies born to smoking and non-smoking mothers are different. HJR Answer: H0 : Average weight of baby from non-smoking = average weight of baby from smoking. Ha : Average weight of baby from non-smoking not = average weight of baby from smoking.

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
```
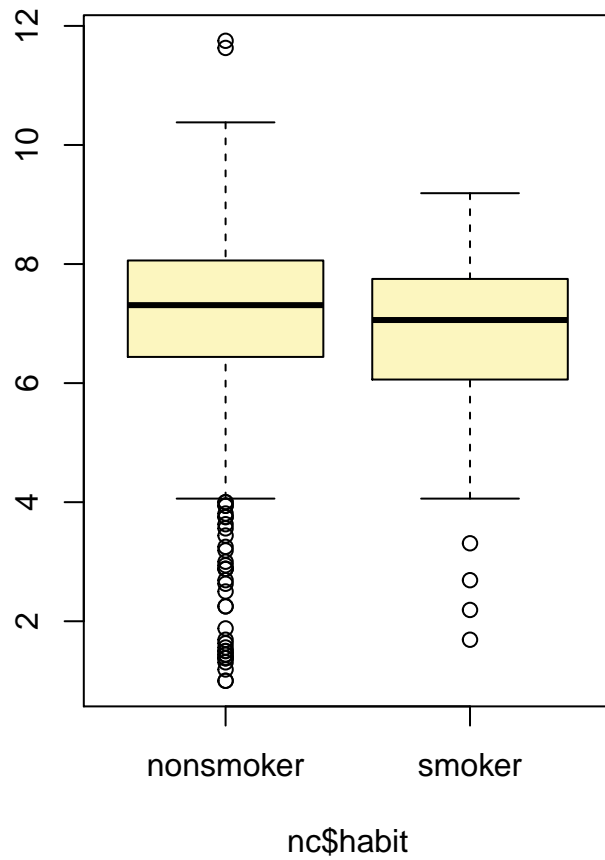
nc$habit

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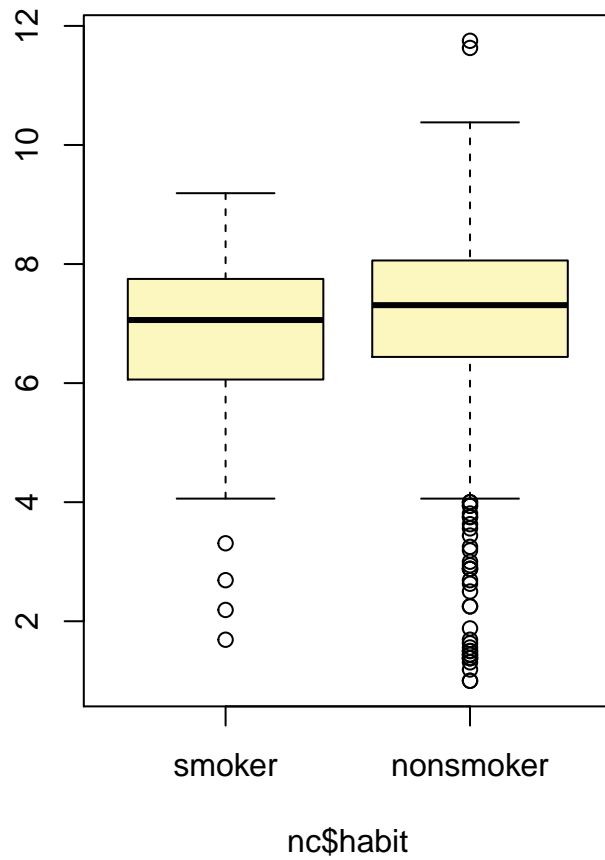
```
## Observed difference between means (nonsmoker-smoker) = 0.3155
##
## Standard error = 0.1338
## 95 % Confidence interval = ( 0.0534 , 0.5777 )
```

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
```

nc$habit

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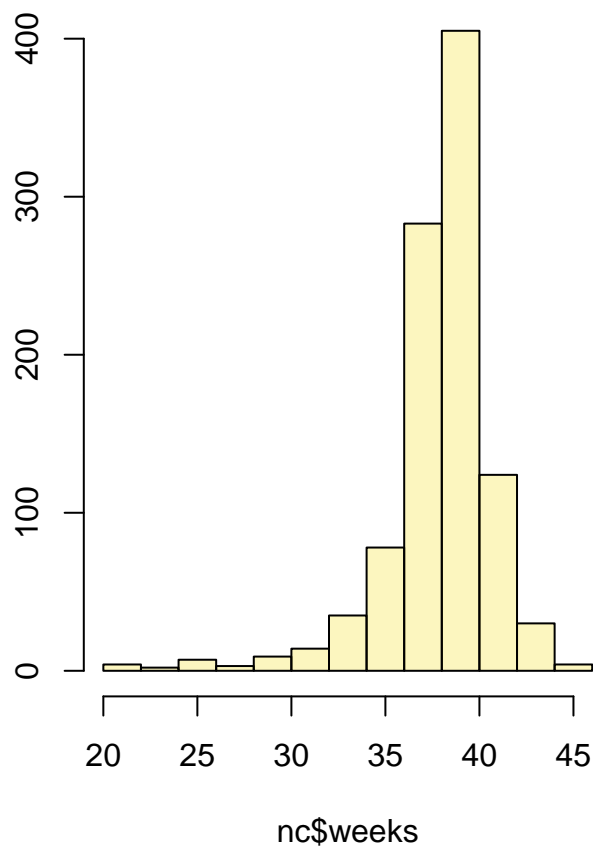
---

**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. JR Answer:

95 % Confidence interval = ( 38.1528 , 38.5165 )

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

```
## Single mean
## Summary statistics:
```
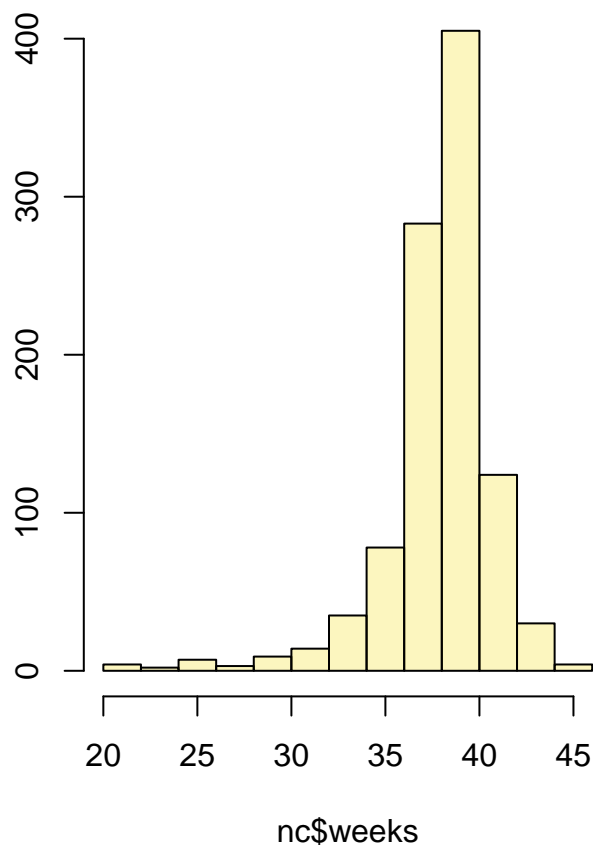
nc$weeks

```
## mean = 38.3347 ;  sd = 2.9316 ;  n = 998
## Standard error = 0.0928
## 95 % Confidence interval = ( 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`.

JR Answer: 90 % Confidence interval = ( 38.182 , 38.4873 )

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

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

nc$weeks

```
## 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.
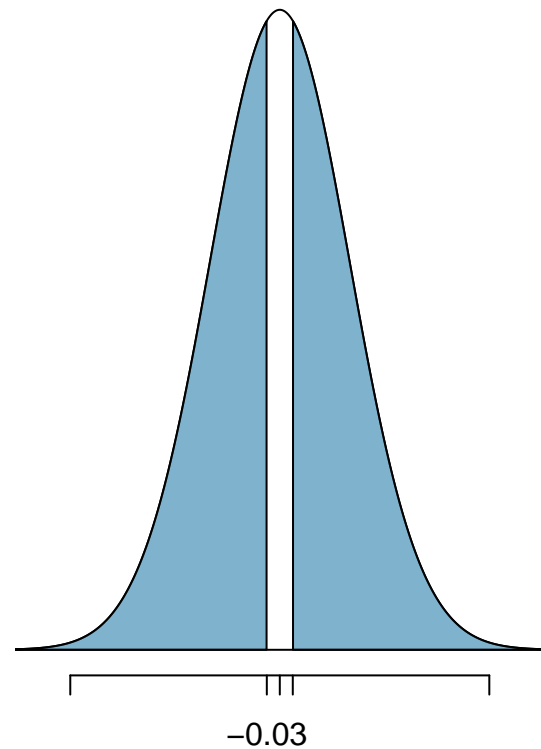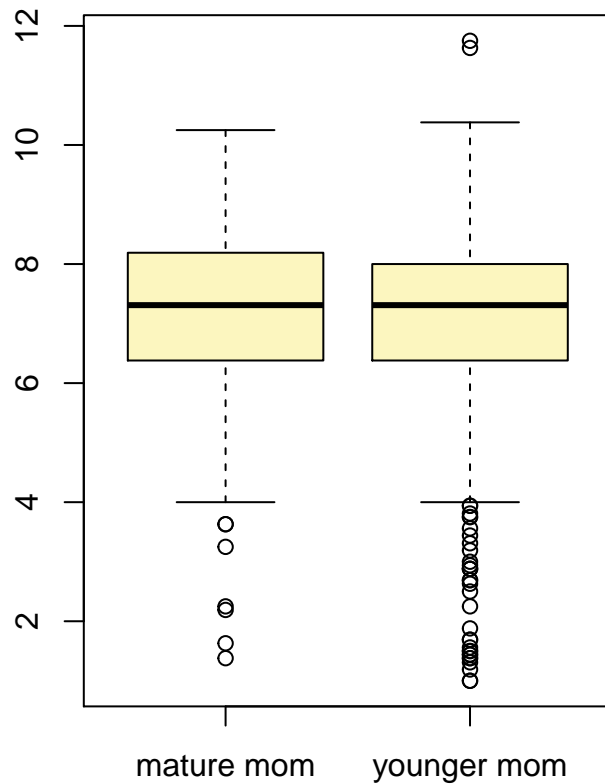
JR Answer: H0: mu_mature mom - mu_younger mom = 0 HA: mu_mature mom - mu_younger mom != 0

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

```
## 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
```

nc$mature

p-value is 0.8526. So, .8526 > .05 so we fail to reject the null hypothesis which was that there was no difference in the weights

- 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. JR Answer: The cutoff is 34 fro younger mother.

```
older <- subset(nc, mature == "mature mom")
younger <- subset(nc, mature == "younger mom")
summary(older$mage)
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   35.00   35.00   37.00   37.18   38.00   50.00
```

```
summary(younger$mage)
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   13.00   21.00   25.00   25.44   30.00   34.00
```

- 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.

```
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
##
##
##
##
```

JR Answer: H0: mu_fage from a premie - mu_fage from a full term = 0 HA: mu_fage from a premie -
mu_fage from a full term != 0

```
inference(y = nc$fage, x = nc$premie, est = "mean", type = "ht", null = 0, alternative = "twosided", met
```

```
## Response variable: numerical, Explanatory variable: categorical
## Difference between two means
## Summary statistics:
## n_full term = 714, mean_full term = 30.2423, sd_full term = 6.6329
## n_premie = 114, mean_premie = 30.3158, sd_premie = 7.5859

## Observed difference between means (full term-premie) = -0.0735
##
## H0: mu_full term - mu_premie = 0
## HA: mu_full term - mu_premie != 0
## Standard error = 0.753
## Test statistic: Z =  -0.098
## p-value =  0.9222
```

nc$premie