

Live data analysis demonstration

Introduction

I thought it would be good to get hands on as early as possible, and to do so for something directly relevant to yourselves. So, we're going to attempt something quite ambitious – in the next two hours or so, we'll go through a whole data analysis from start to end. Lets get started.

Meta-task

Write things that you don't understand, and need to know about in subsequent classes. We will cover these things.

The question

What should be our question? As always, there are some influences and some constraints. We should ask a question of interest to us, and of some importance. And we should be able to collect the data, within our current constraints, necessary to answer the question.

The question we will address is “do male and female reaction times of students at the University of Zurich differ?”.

Why this question? Reaction times are important, safety, sport...

Expectation

We can have a look on the internet, and pretty easily find lots of studies of reaction times and gender (e.g., A comparative study of visual and auditory reaction times on the basis of gender and physical activity levels of medical first year students (by the way, we will later in the course critique this paper – it has some pretty poor features).

Generally, we see that males tend to have faster reaction times than females. So we expect that to be the same for students at the University of Zurich.

Given that you know this pattern, and you are the subjects, its interesting to see if you women can buck the trend, perhaps by trying especially hard. Though know the men know you might do this, it probably won't work!

How are we going to present the results?

Thinking backwards from how we present the final results, can often be quite useful.

I think a nice box and whisker plot will work here (Owen will sketch this). We will have two groups of reaction times. Put another way, we will have one explanatory variable (gender) and one response variable (reaction time). The explanatory variable gender is a categorical / discrete variable. The response variable is continuous.

We expect the distribution of male reaction times to be have a lower mean than the distribution of female reaction times.

We will look at this graph, and answer our question (wow, without any statistical test – YES!).

What statistical test will we use?

Reaction times (the response variable) we expect to be quite normally distributed, though cannot be negative. Gender will be categorical with two levels (male and female). We don't expect greater or less variation in reaction times among males compared to among females.

Based on these expectations, we will use a linear model, which assumes normally distributed residuals and equal variances among groups. The traditional name for the test is the T-test.

Based on convention, and little else, we will say there is a significant difference between male and female reaction times if the observed difference has a p-value of less than 0.05.

Selection of subjects

We usually need to very carefully select the subjects of our study. Ideally, as we're interested in students at the University of Zurich (see the question), we would select a cross section of such students. Instead, you are going to be the subjects, and you are not representative of all students. You're relatively young, on average, you're studying natural sciences, etc. So we will have to be very cautious if we make statements about students at the University in general, and perhaps one might even now conclude that we can't really answer the question.

Perhaps we need the question "do male and female reaction times of biology and biomedicine students, in their first year, at the University of Zurich differ?"

Ethical clearance and considerations

If we aimed to publish these results, or in some other way disseminate them, we would need ethical clearance for research involving humans. We're not, so we don't have ethical clearance.

However, please do not include any personal information in any of the data you contribute to the exercise.

Data collection

Create for yourself a unique ID code, so that if we want to collect more data about you, we could related the reaction time data to that. Write this code down somewhere safe, keep it.

Go to the Human Benchmark website.

Do the reaction time test and write down your reaction time in fractions of a second.

While you're there, also please do the other three tests, and write down your score. We don't need these for our current question about reaction times, but we might look at this data later in the course.

Now go to this web page: <https://goo.gl/forms/cIIdMtTj1g1npfLk2> and enter your ID, gender, and four scores. Please be careful!

Go ahead and do all that.

Look at the data!

Here is the link to the datasheet containing all the data you just recorded: <https://docs.google.com/spreadsheets/d/1wltFGU6bbkYeiDzzvQTIsYwOupqDZL8xiLAX0QJWWk/edit#gid=1188775314>

Lets have a look at it, see what you've done. (I'm scared! I know how difficult it is to enter data without making mistakes! And I have some experience of how different people are!)

Lets get the data into our data analysis software of choice (R, via RStudio)

First note that Owen has gone to the responses googlesheet, and in the “File” menu, clicked on “Publish to web...” and chosen to publish as “Comma separate values”.

Now we can read that web page of comma separated values into R:

```
## First we load a required package (we need to install this if we haven't already)
library(readr)
## also we will use some other packages
library(dplyr)
library(ggplot2)

## Now read in the data, using the read_csv() function. We give it the URL of the published version of
the_URL <- "https://docs.google.com/spreadsheets/d/e/2PACX-1vQFgYX1QhF9-UXep22XmPow1ZK5nbFHix9nkQIa0DzqI
class_RT<sup>s</sup> <- read_csv(the_URL)

## Have a look at the data in R
#View(class_RT<sup>s</sup>)
## or just do
class_RT<sup>s</sup>

## # A tibble: 8 × 7
##       Timestamp `Please enter the unique ID code you gave yourself.`
##       <chr>                                <chr>
## 1 24/11/2016 14:30:56 asdfasdf
## 2 24/11/2016 15:02:36 asdfasdasdf
## 3 24/11/2016 15:02:36 a
## 4 24/11/2016 15:02:36 b
## 5 24/11/2016 15:02:36 c
## 6 24/11/2016 15:02:36 d
## 7 24/11/2016 15:02:36 e
## 8 24/11/2016 15:02:36 f
## # ... with 5 more variables: `What is your gender?` <chr>, `Please enter
## #   your average reaction time in seconds (e.g., 0.326).` <dbl>, `Please
## #   enter your score on the Verbal Memory test.` <int>, `Please enter your
## #   score on the Number Memory test` <int>, `Please enter your score on
## #   the Visual Memory test.` <int>
```

Now we need to do some data wrangling (cleaning and tidying)

Clean up the column / variable names:

```
## Must be very careful to get the next line right!!! Really important!!!
names(class_RT<sup>s</sup>) <- c("Timestamp", "ID", "Gender", "Reaction_time",
                          "Verbal_memory_score", "Number_memory_score",
                          "Visual_memory_score")
class_RT<sup>s</sup>

## # A tibble: 8 × 7
##       Timestamp      ID Gender Reaction_time Verbal_memory_score
##       <chr>      <chr> <chr>      <dbl>          <int>
## 1 24/11/2016 14:30:56 asdfasdf  Male         0.45            42
## 2 24/11/2016 15:02:36 asdfasdasdf Male         0.55            10
```

```
## 3 24/11/2016 15:02:36      a Female      0.23      42
## 4 24/11/2016 15:02:36      b  Male      0.10      42
## 5 24/11/2016 15:02:36      c Female      0.30      42
## 6 24/11/2016 15:02:36      d  Male      0.32      42
## 7 24/11/2016 15:02:36      e Female      0.45      10
## 8 24/11/2016 15:02:36      f Female      0.29      42
## # ... with 2 more variables: Number_memory_score <int>,
## #   Visual_memory_score <int>
```

Check the variable types are correct.

- Timestamp should be a character
- ID should be a character
- Gender should be a character
- The remaining four variables should be numeric (if fractional, if whole numbers).

```
class_RTs
```

```
## # A tibble: 8 × 7
##       Timestamp      ID Gender Reaction_time Verbal_memory_score
##       <chr>      <chr> <chr>      <dbl>          <int>
## 1 24/11/2016 14:30:56 asdfasdf  Male      0.45           42
## 2 24/11/2016 15:02:36 asdfasdasdf Male      0.55           10
## 3 24/11/2016 15:02:36      a Female      0.23           42
## 4 24/11/2016 15:02:36      b  Male      0.10           42
## 5 24/11/2016 15:02:36      c Female      0.30           42
## 6 24/11/2016 15:02:36      d  Male      0.32           42
## 7 24/11/2016 15:02:36      e Female      0.45           10
## 8 24/11/2016 15:02:36      f Female      0.29           42
## # ... with 2 more variables: Number_memory_score <int>,
## #   Visual_memory_score <int>
```

Correct or exclude problematic data

If we have problems here, with variables of the wrong type, it probably means some of the data entry is a bit messed up.

```
## Have to do this live!!!
## e.g. to exclude observations with character entries in Reaction_time variable
class_RTs <- filter(class_RTs, !is.na(as.numeric(Reaction_time)))
```

Once fixed, we need to make the variable have the correct type

```
## try using type_convert() from readr package.
class_RTs <- type_convert(class_RTs)
```

Check the number of observations

Should be the same as we saw in the datasheet, which should be number of you in this classroom.

The number of observations and variables is given by R in the first line of output when we type the name of the data object:

```
class_RTs
```

```
## # A tibble: 8 × 7
##       Timestamp      ID Gender Reaction_time Verbal_memory_score
##       <chr>      <chr> <chr>      <dbl>          <int>
## 1 24/11/2016 14:30:56 asdfasdf Male         0.45           42
## 2 24/11/2016 15:02:36 asdfasdasdf Male         0.55           10
## 3 24/11/2016 15:02:36      a Female         0.23           42
## 4 24/11/2016 15:02:36      b Male          0.10           42
## 5 24/11/2016 15:02:36      c Female         0.30           42
## 6 24/11/2016 15:02:36      d Male          0.32           42
## 7 24/11/2016 15:02:36      e Female         0.45           10
## 8 24/11/2016 15:02:36      f Female         0.29           42
## # ... with 2 more variables: Number_memory_score <int>,
## #   Visual_memory_score <int>
```

Visualise the data

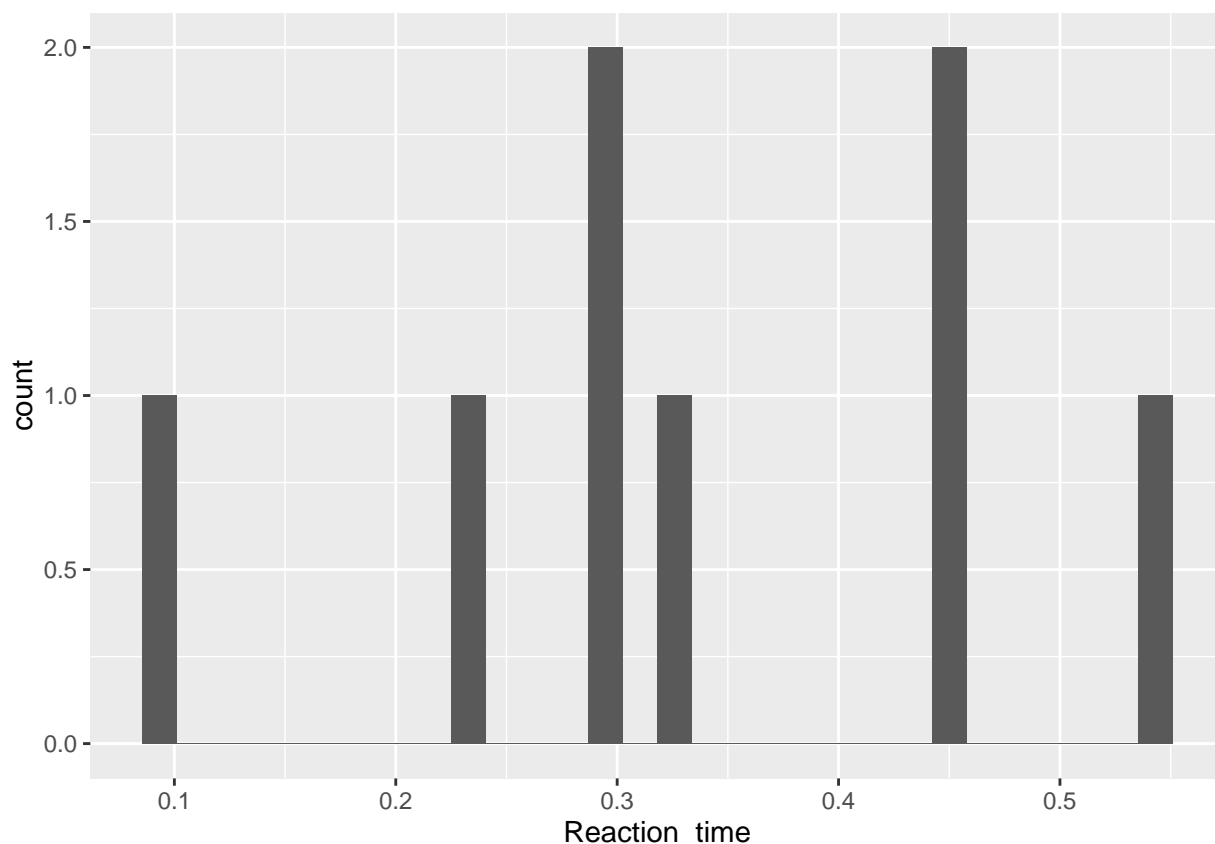
When we visualise the data, we're trying to do at least three things, and are not trying to do at least one.

We're not trying to make the most beautiful graph in the world, so we can put it in our report / presentation etc. We just want to clearly see the data.

We are trying to 1) do further checks for possible errors in the data, 2) making some initial assessments of how the data is distributed, 3) see what we think is the answer to our question.

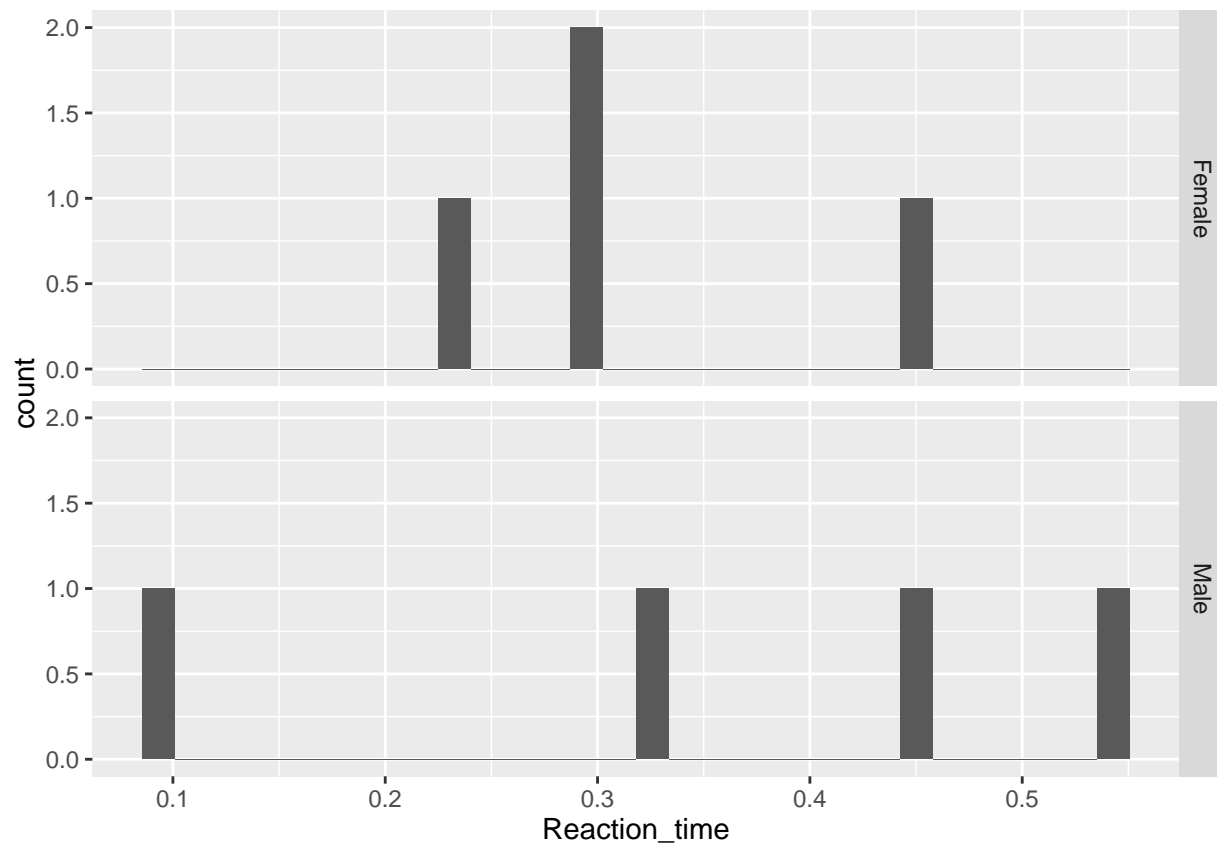
A histogram of all the data:

```
qplot(x=Reaction_time, data=class_RTs)
```



Separate histograms for each gender:

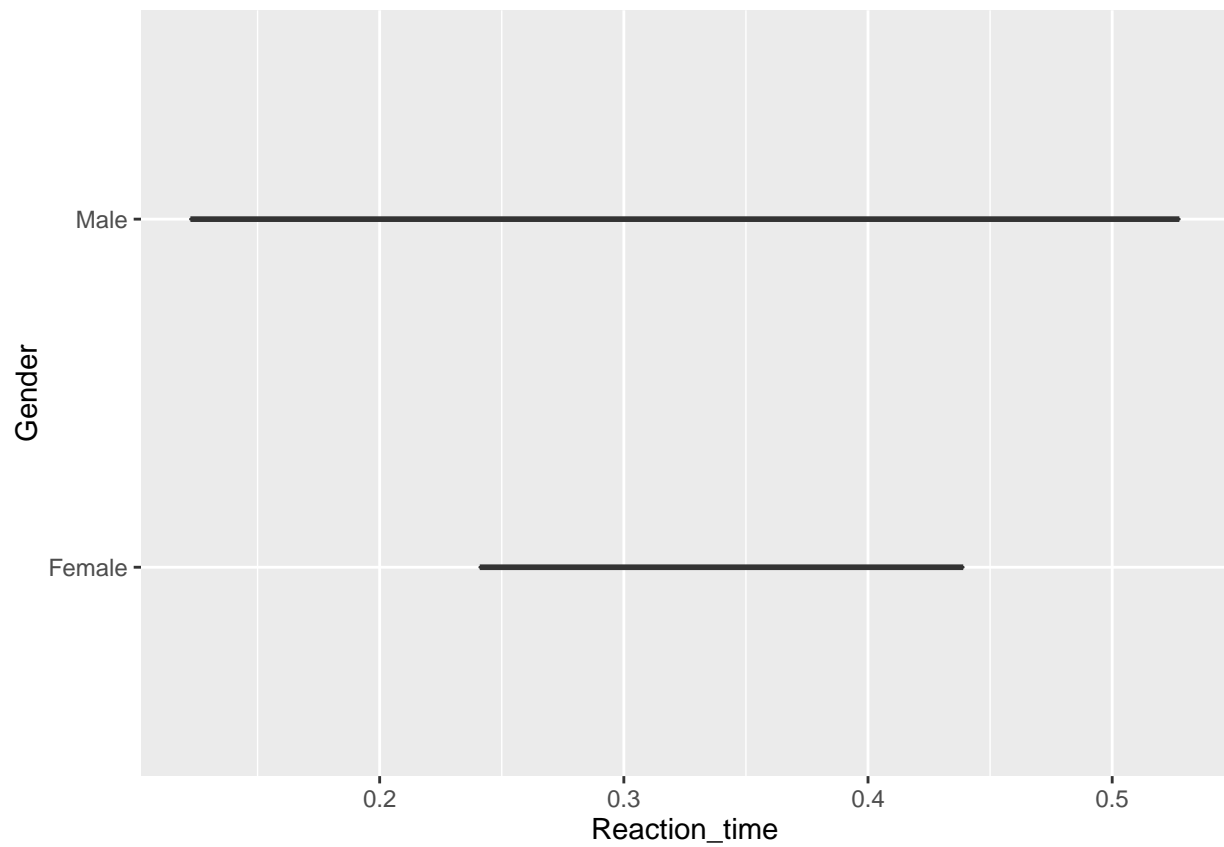
```
qplot(x=Reaction_time, data=class_RTs, facets = Gender ~ .)
```



A box and whisker plot:

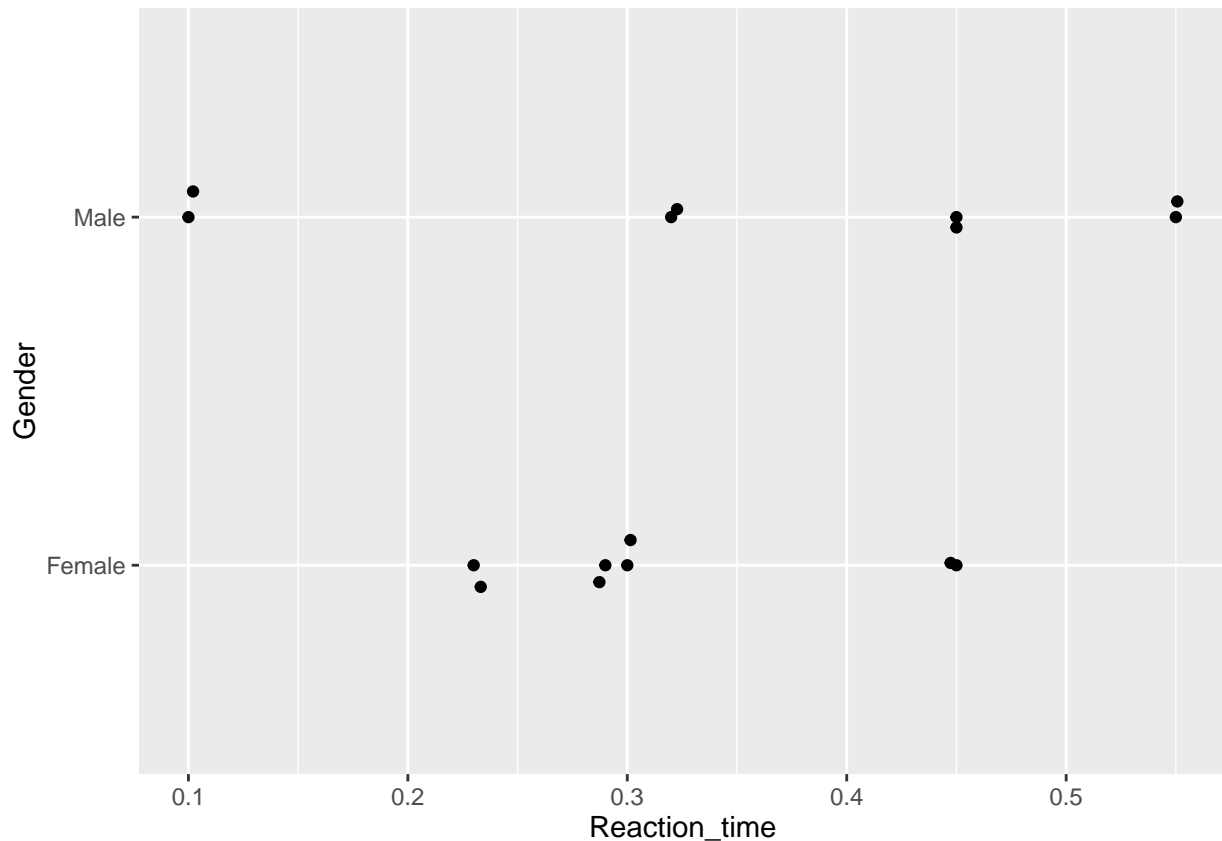
```
qplot(x=Reaction_time, y=Gender, data=class_RTs, geom = "boxplot")
```

```
## Warning: position_dodge requires non-overlapping x intervals
```



Or just the data points (with some jitter, to separate overlapping points):

```
qplot(x=Reaction_time, y=Gender, data=class_RTs) + geom_jitter(height=0.2)
```



What do we think about the three things? Any likely errors? How is the data distributed (within and between groups)? Does it look like there is a difference in reaction times (if so, by how much on average, and which group is faster)?

Assess assumptions

Before we even start to think about running a statistical test, we must check if the specific test we intend to run is justified. That is, we must check if the assumptions of the test are likely to be met.

Independence

The t-test assumes that observations are independent.

How was the data collected (hopefully not one gender on one day, and the other on another day, or something similarly confounding)?

Do we have more than one observation per subject? Not in this case, because you typed in the average. But we could have. Then the observations from the same individual would not be independent of each other. They would share in common the person they originated from. This would make the statistical test unreliable.

Normally distributed residuals

We can get a good idea about this, in this case, by looking at the distribution of the two groups of reaction times (see above). Obviously we need to have in mind some idea of what the normal distribution looks like, and how close the data have to look like one. There are quantitative tests for normality, we may look at them later.

Equal variance

The spread of the reaction times for men, and the spread for women, should be about the same.

We can get a good idea about this, in this case, by looking at the distribution of the two groups of reaction times (see above). Again, we need to have in mind how similar the variance can be, without invalidating this assumption the data. There are quantitative tests for equal variance, we may look at them later.

Do the statistical test

We have to do a test, or more generally, some statistics, to give some kind of assessment of certainty / uncertainty in our answer. Traditionally, this is done with a p-value, and if its lower than 0.05 we say the result is significant (i.e. the results are very consistent with no difference). If its higher than 0.05 we accept the null hypothesis that there is no difference.

Another way to quantify uncertainty, is to give the difference in the means of the two groups, and a measure of certainty in this difference. If the difference between the means close to zero, and the uncertainty overlaps zero, then we conclude there is no strong difference.

We'll do this with a T-test, as we already planned. Before we go on, there is something very important we should figure out, and we should do this every time before we run a statistical test. Figure out the degrees of freedom.

There will be learning about this later in the course. For now, know that for a t-test the degrees of freedom are the number of observations minus two. Here that is $8 - 2 (= 6)$. This is really important to figure out in advance, as its a great way to check that R is doing the test we think we're telling it to do.

```
my_ttest <- t.test(Reaction_time ~ Gender, data=class_RTs, var.equal=TRUE)
my_ttest
```

```
##
## Two Sample t-test
##
## data: Reaction_time by Gender
## t = -0.34771, df = 6, p-value = 0.7399
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.3013952 0.2263952
## sample estimates:
## mean in group Female mean in group Male
## 0.3175 0.3550
```

Lots of information there. We will teach you how to read this in later lectures.

For now, we can find the p-value: 0.73993. And the difference between the means: 0.038 and the lower (-0.301) and upper (0.226) 95% confidence limits on that difference.

Report and communicate the results

The results as a sentence

We should write a sentence that gives the direction and extent of difference, and a measure of certainty / uncertainty in that finding. It is totally unacceptable, though common, to just write “there was a significant difference”. If we want to give a p-value (and most people tend to expect to see one), we should remind about the statistical test used (remind because we may have already mentioned it) and give the degrees of freedom, the value of the test statistic, and the p-value.

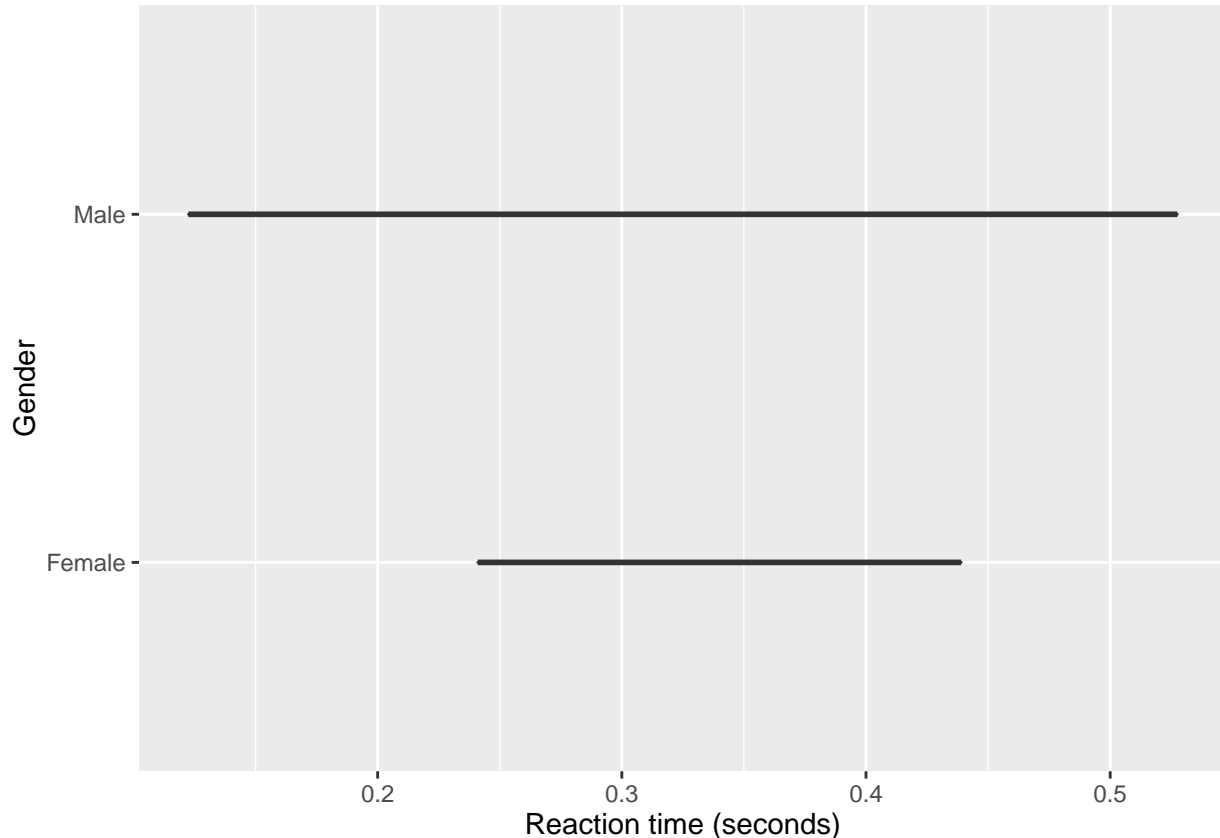
Insert sentence here, once we have the results.

The results graphically

The aim here is to make a beautiful graph that very clearly communicates the findings! This doesn't mean "fancy" and or "complex". Often simpler is better. Getting the basic right is essential, of course.

```
qplot(x=Reaction_time, y=Gender, data=class_RT, geom = "boxplot") +  
  xlab("Reaction time (seconds)")
```

```
## Warning: position_dodge requires non-overlapping x intervals
```



Wow! That was easy.

Do not use a table

Here, a table is not necessary. The results are in the sentence and in the graph. Enough already!

Critical thinking

- How might the work be flawed?
- How might the analysis be flawed (assumptions violated)?
- Is the difference (i.e. effect size) small, medium, large, relative to differences caused by other factors?
- How general might be the finding?
- How do the qualitative and quantitative findings compare to those in previous studies?
- What could have been done better?

- What are the implications of the findings?