

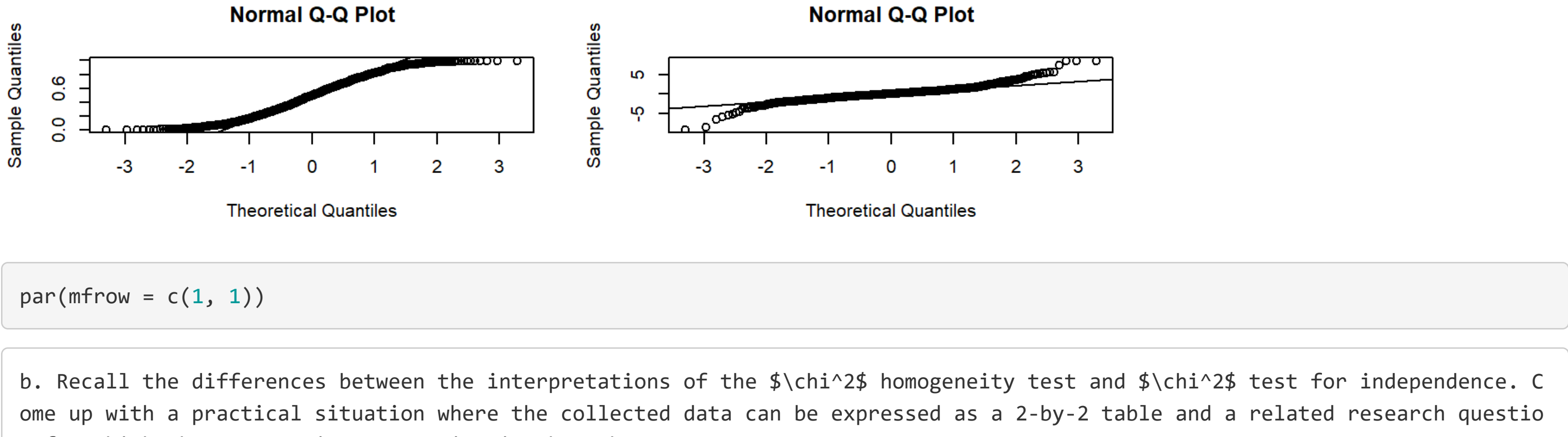
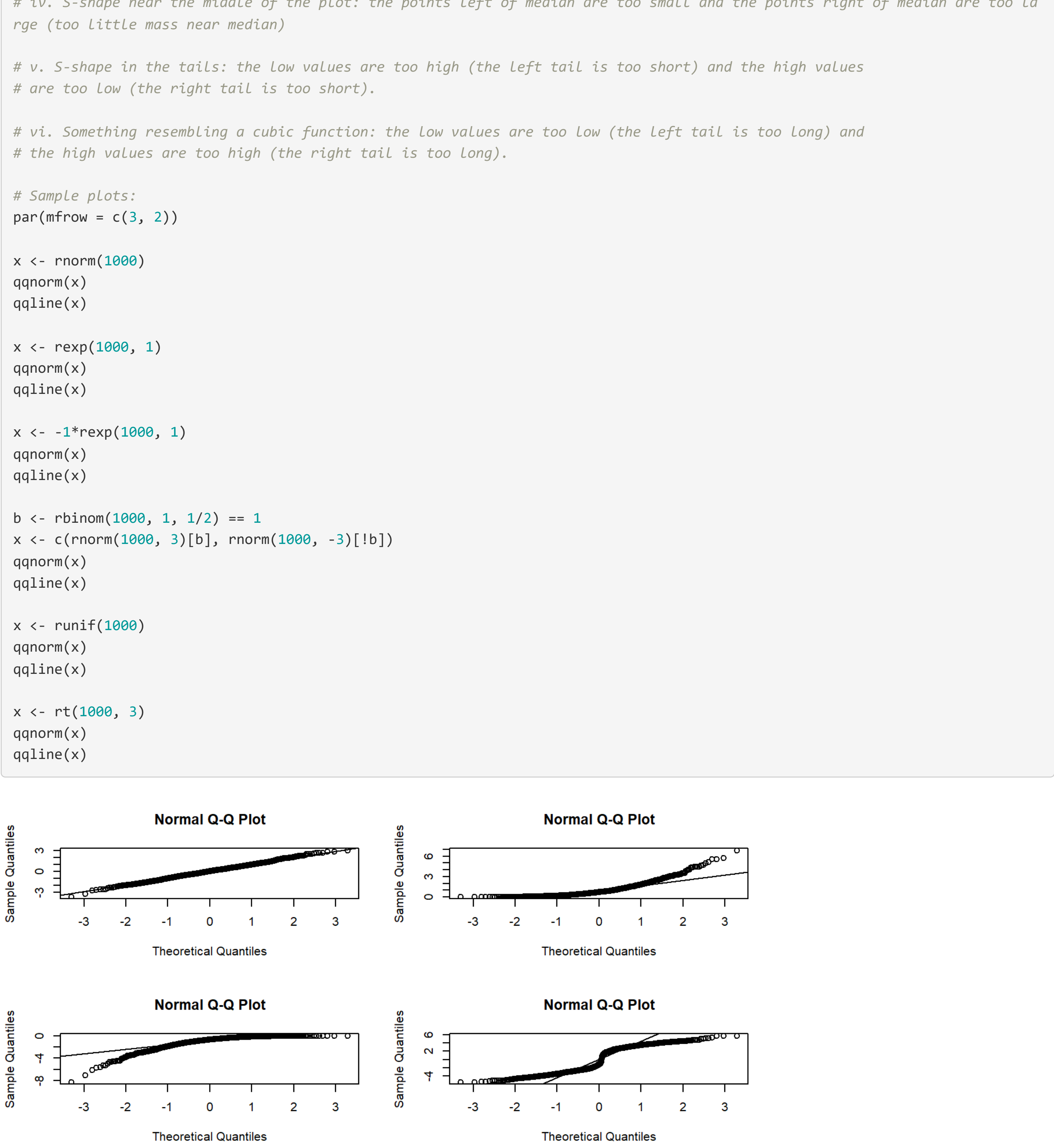
MS-C1620 Statistical Inference

Exercise 6

Homework exercise

To be solved at home before the exercise session.

1. a. Assume that we have an iid. random sample x_1, \dots, x_{1000} and we'd like to use the normal Q-Q plot to assess whether the sample came from a normal distribution. How do you expect the normal Q-Q plot to roughly look like (i.e. what general features do you expect it to have and why), if the true distribution of the data is
- a normal distribution,
 - a right-skew distribution,
 - a left-skew distribution,
 - a bimodal distribution,
 - a distribution with light tails,
 - a distribution with heavy tails?



- b. Recall the differences between the interpretations of the χ^2 homogeneity test and χ^2 test for independence. Come up with a practical situation where the collected data can be expressed as a 2-by-2 table and a related research question for which the correct interpretation is through
- the χ^2 homogeneity test,
 - the χ^2 test for independence.

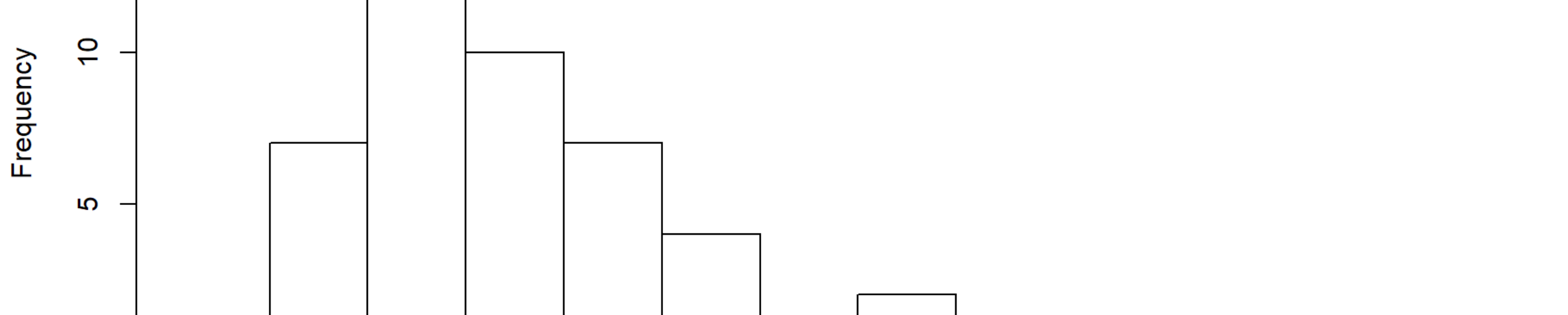
```
# The key difference between the two tests is in how the data is sampled, i.e., are the margins fixed or not.
# E.g. assume we're interested in studying whether sex (female/male) has an effect on the voting preference
# (democrat/republican) in the US and for this we interview n people in the street.
# These data can be collected into a two-by-two table such that the row variable is sex and the column
# variable is voting preference.
# i. If we choose beforehand that we will interview n1 females and n2 males, then studying the independence of the two variables will be questionable (since sex is not fully random anymore with its marginal frequencies fixed). The correct interpretation is through the homogeneity test which compares two populations, in this case female and male, in their voting preferences.
# ii. If we do not choose beforehand the marginal numbers of females and males, sex is a random variable and we can measure its independence with the voting behavior. The correct interpretation is now through the test for independence.
```

Class exercise

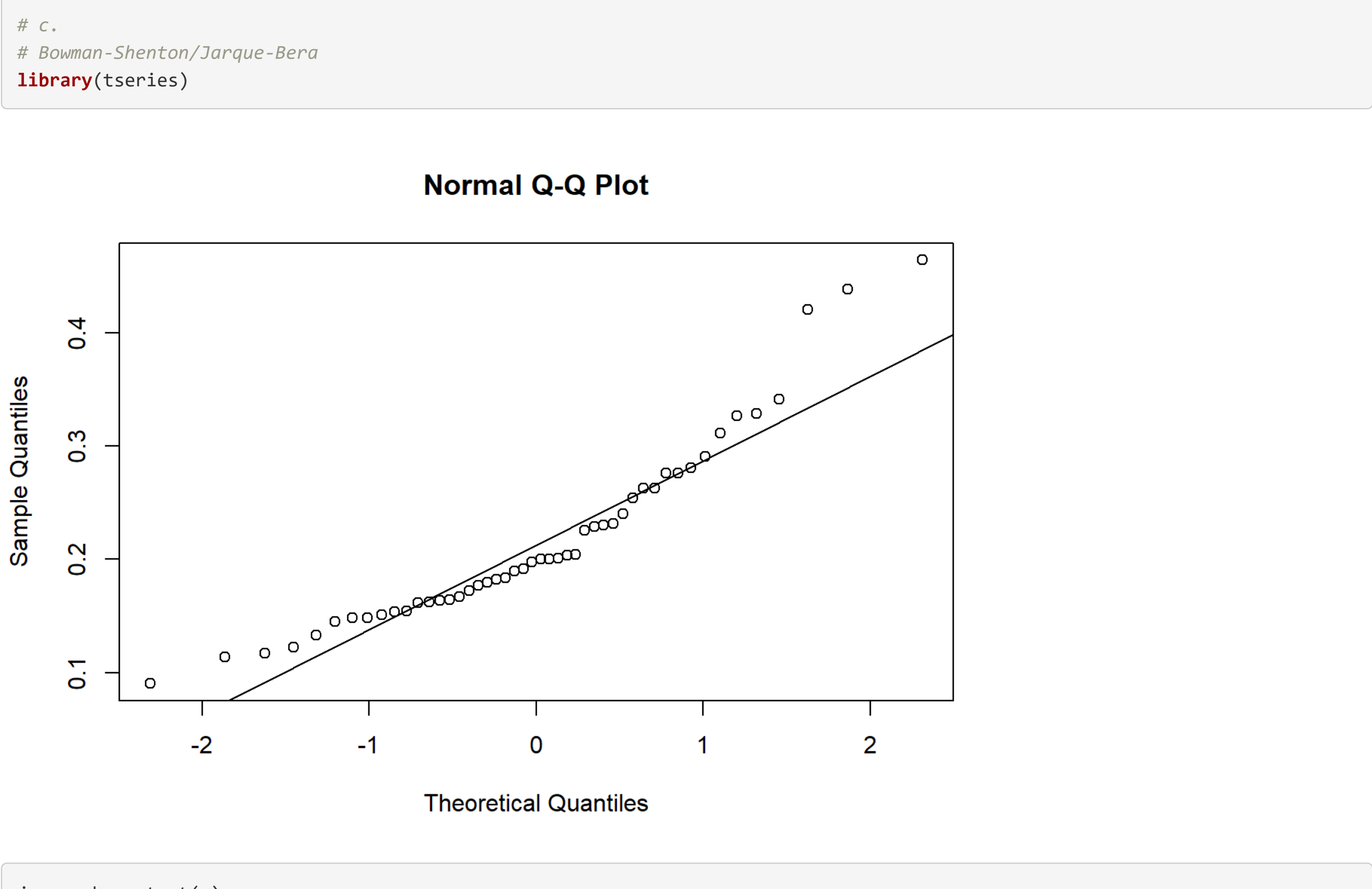
To be solved at the exercise session.

Note: all the needed data sets are either given below or available in base R.

1. The data set `rock` contains measurements on 48 rock samples from a petroleum reservoir. Treat the data as an iid. random sample from some distribution and test whether the distribution of shape is normal.
- Visualize the data to obtain a preliminary idea of the possible normality of the data.
 - Use the normal Q-Q plot to gain more evidence on the normality/non-normality of the data.
 - Conduct the Bowman-Shenton (Jarque-Bera) and the Shapiro-Wilk tests of normality on significance level 0.05.
 - After all the previous, would you conclude the data to be normal (or normal enough for methods with normality assumptions)?
 - Why is the data not really iid.?



- b. The normal Q-Q plot gives more evidence of positive skewness (it is reminiscent of part ii in Homework 1a.)
- ```
qqnorm(x)
qqline(x)
```
- c. Bowman-Shenton/Jarque-Bera
- ```
library(tseries)
```



```
jarque.bera.test(x)
```

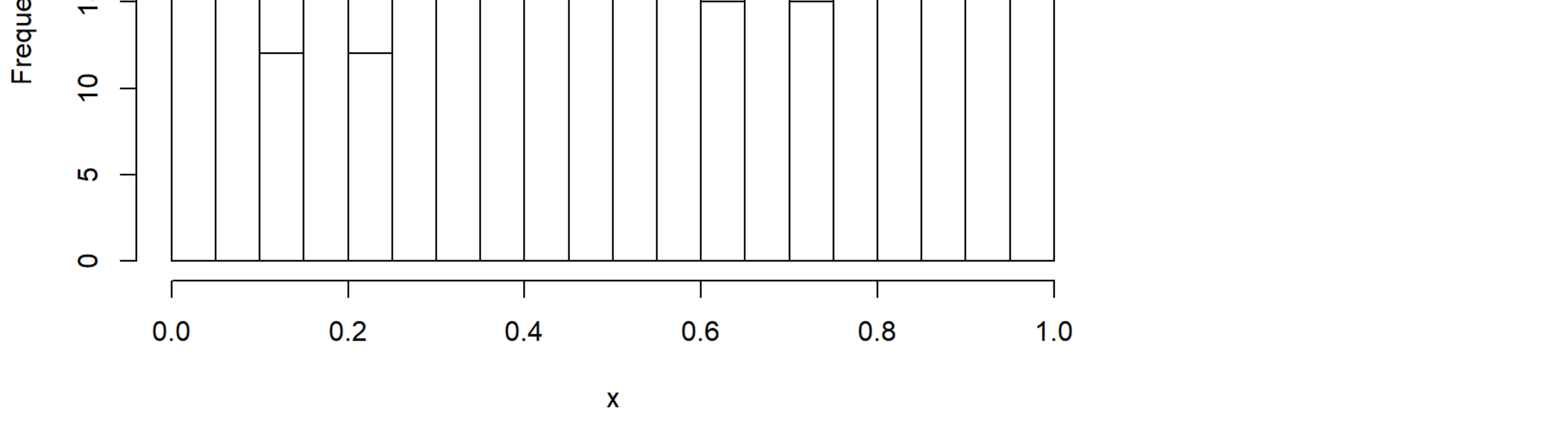
```
##
## Jarque Bera Test
##
## data: x
## X-squared = 13.402, df = 2, p-value = 0.00123
```

```
# Shapiro-Wilk
shapiro.test(x)
```

```
##
## Shapiro-Wilk normality test
##
## data: x
## W = 0.90407, p-value = 0.000531
```

- d. Both tests in c. reject their null hypotheses of normality. Based on all the previous evidence, the data can not be deemed normal enough to rely on normality assumptions in any further analyses.
- (Note that it is a different matter whether the next analysis steps involve methods that allow the normality assumption to be "covered" by large enough sample size (by the central limit theorem).)
- e. See the help file of the dataset: The sample is not iid. as, to obtain the 48 measurements, first 12 "core samples" were obtained (randomly?) and then from each of these 4 observations were taken to yield the final 48 observations. Thus, the sets of 4 observations come from a same core sample and as such are not independent, even if the different core samples were.

2. The data set `randu` contains 400 triples of successive random numbers from the random number generator `RANDU`. Use the χ^2 goodness-of-fit test to assess whether the first elements in the triplets really obey the uniform distribution on $[0, 1]$.
- Extract the first elements in the triplets and visualize their sample distribution.
 - Discretize the values into a suitable number of categories and calculate the observed category frequencies.
 - Compute the corresponding expected category probabilities under the uniform distribution on $[0, 1]$.
 - Recall the hypotheses of the test and conduct it on significance level 0.05.
 - What are the conclusions of the test? Compare your results with someone who used a different choice of categories for the discretization.



```
# b.
# The categories: [0, 0.1], (0.1, 0.2], ...
obs_x <- diff(sapply(seq(0, 1, by = 0.1), function(i) sum(x <= i)))

# c.
# Each category has the same width 0.1 and the postulated distribution is uniform on [0, 1]
# -> the expected category probabilities are each 0.10 (40 obs. per category)
p_x <- rep(0.1, 10)

# d.
# The null hypothesis is that the data was generated by the uniform distribution on [0, 1]
# and the alternative that it was not.
chisq.test(x = obs_x, p = p_x)
```

```
##
## Chi-squared test for given probabilities
##
## data: obs_x
## X-squared = 7.85, df = 9, p-value = 0.5493
```

- e. The test p-value is 0.5493
- > no evidence against H_0 , it is still plausible that the data is from $Uniform[0, 1]$.
- By choosing the categories suitably, it is most likely possible to get the opposite result (recall the Type I and II errors). However, this should not be taken advantage of in practice...

3. The data set `Titanic` contains information on the fate of passengers on the fatal maiden voyage of the ocean liner "Titanic". We use the data to study whether there is a connection between the sex (Male/Female) of a passenger and surviving from the ship (No/Yes).
- Extract a marginal table containing only the cross-tabulation of the variables `Sex` and `Survived`.
 - Find a suitable way to visualize the data.
 - Which test is appropriate for these data (and why?). χ^2 homogeneity test or the χ^2 test for independence?
 - Conduct your chosen test on significance level 0.05 and state your conclusions.



```
# c.
# It seems plausible that there were no quotas on Female/Male passengers on the ship. As such, both factors had
# their margins non-fixed and the correct test is the test for independence.

# d.
chisq.test(x)
```

```
##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data: x
## X-squared = 454.5, df = 1, p-value < 2.2e-16
```

- Very Low p-value
- > sex and survival status are not independent -> females had stat. significant higher chance of surviving.
4. (Optional) Choose your favorite non-normal distribution and use simulations to study the Type II error probabilities of the Bowman-Shenton (Jarque-Bera) and Shapiro-Wilk tests of normality for that distribution on different sample sizes (e.g. $n = 10, 100, 1000, 10000$). That is, find out the probability of falsely concluding that the data comes from a normal distribution when it does not.