

Stat 5810, Section 003
Statistical Visualization I
Fall 2018
Homework 1

ShaunMicheal Bartschi

A01975136

October 17, 2018

Homework Assignment 1 (10/6/2018)

60 Points — Due Wednesday 10/17/2018 (via Canvas by 11:59pm)

- (i) (48 Points) In this question, you have to work with the *Hidalgo1872.rda* data set that is provided in Canvas. You can also directly download the data from <https://github.com/cran/MMST/tree/master/data>. Do not use any other version of this data set you may find on the internet or your results may differ slightly as different versions of the data set seem to exist. This data set was originally a part of the *MMST* R package, but that package was removed from the CRAN repository.

The data set consists of 485 rows and 3 columns. The first column (*thickness*) contains the thickness of 485 Mexican stamps, measured to a precision of a thousandth of a millimeter. It is my understanding that column 2 (*thicknessA*) repeats these measurements if the stamp was pressed in 1872 (and contains NA otherwise). Column 3 (*thicknessB*) repeats these measurements if the stamp was pressed in 1873/74 (and contains NA otherwise). There are a few data entry errors where both columns 2 and 3 contain a measurement.

The full (!?) description of the data set, as originally provided in the *MMST* R package, is shown in Figure 1.

- (a) (1 Point) Load the *Hidalgo1872* data set and all required R packages to answer this question. Show your R code.

```
> setwd("C:/Users/Shawn/Desktop/StatVis/HW1")
> load("Hidalgo1872.rda")
> library(ggplot2)
```

- (b) (2 Points) Draw a basic histogram for *thickness* using *ggplot2*. Include your R code and the resulting graph.
- (c) (6 Points) Further improve your histogram from (b). Try three different binwidths: 0.001, 0.002, and 0.005 and use 0.0005, 0.001, and 0.0025 as the center, respectively. You should also adjust the range of the horizontal and vertical axes, labels, title, etc. **Clearly indicate which changes you made and why you made these changes.** As in class, make these changes

Hidalgo1872	MMST HIDALGO 1872 STAMP DATA
Description	
Hidalgo postage stamps, 93, 96, 98	
Usage	
<code>data(Hidalgo1872)</code>	
Format	
A data frame with 485 observations on the following 3 variables.	
<code>thickness</code> a numeric vector	
<code>thicknessA</code> a numeric vector	
<code>thicknessB</code> a numeric vector	
References	
A. Izenman (2008), <i>Modern Multivariate Statistical Techniques</i> , Springer	
Izenman, A.J. and Sommer, C.J. (1988). Philatelic mixtures and multimodal densities, <i>Journal of the American Statistical Association</i> , 83 , 941-953.	

Figure 1: Description of the Hidalgo 1872 data set, obtained from <https://mran.microsoft.com/snapshot/2014-09-30/web/packages/MMST/MMST.pdf>.

step by step and continue with further adjustments until your graph is ready for publication. Include your three final graphs for these three binwidths and the R code for these final graph. No need to include any intermediate graphs and the R code for those. Hint: When you get warnings from ggplot2, check carefully. It is easy to cut off parts of the histogram on the horizontal or vertical axis, in particular when you copy and paste your R code and forget to adjust some of the arguments.

- (d) (1 Point) Repeat (b) from above, now using the *hist* function from baseR.
- (e) (6 Points) Repeat (c) from above, now using the *hist* function from baseR.
- (f) (4 Points) Based on your final histograms in (c) and (e), how many modes does the Hidalgo 1872 data set seem to have? Answer this question separately for your three different binwidths. What is your overall conclusion regarding the number of modes?
- (g) (4 Points) Recall that the stamps originate from the years 1872 and 1873/74. The help page shown in Figure 1 is not very helpful, so we have to make our own assumptions: If the third column (*thicknessB*) contains a value, then this is a measurement from 1873/74. Otherwise, it is a measurement from 1872. Using any approach in R you are familiar with, add a column called

Year to the `Hidalgo1872` data frame. You cannot modify the data outside of R, e.g., via Excel. Show your R code, the first 6 lines of your modified data frame, and a table that summarizes the new *Year* column. Hint: There should be 289 measurements for 1872 and 196 for 1873/74.

- (h) (6 Points) Start with a basic histogram for *thickness*, conditioned on the two options for *Year*, using `ggplot2`. Optimize this graph in multiple steps. Use a layout that shows the two resulting histograms above each other for better comparison. As before, try three different binwidths: 0.001, 0.002, and 0.005 and use 0.0005, 0.001, and 0.0025 as the center, respectively. Include your R code and the resulting final figure that consists of six histograms overall: 1872 on top and 1873/74 at the bottom and binwidths 0.001, 0.002, and 0.005 from left to right. The two histograms above each other should follow the small multiple principle, i.e., have the same ranges for the horizontal and vertical axes. The histograms besides each other do not have to follow this principle
- (i) (6 Points) Repeat (h) from above, now using the `hist` function from `baseR`. In particular, use the same layout as described above. Hint: Now you have to be really careful to use the proper ranges for the horizontal and vertical axes.
- (j) (4 Points) Based on your final histograms in (h) and (i), how many modes does the Hidalgo 1872 data set seem to have for each year? Answer this question separately for your two different years (1872 and 1873/74) and your three different binwidths. What is your overall conclusion regarding the number of modes?
- (k) (4 Points) Would boxplots be good replacements for the two histograms for the two years? Compare carefully what can be seen in the box plots and what cannot be seen. First create two basic boxplots with a package of your choice. Then refine them. Add labels as needed and make sure your boxplots follow the small multiple principle. As always, include your R code and the final resulting graphs. Then answer **yes** (they are good replacements) or **no** (they are not good replacements). Justify your answer!
- (l) (4 Points) Have you ever heard of violin plots? If not, google them! Find a suitable R package that creates violin plots or see how they can be created in `ggplot2`. Would violin plots be good replacements for the two histograms? First create two basic violin plots with a package of your choice. As always,

include your R code and the final resulting graphs. Then answer **yes** (they are good replacements) or **no** (they are not good replacements). Justify your answer!

(ii) (12 Points) This question makes use of the *Pima.tr2* data set from the *MASS* R package again. These graphs may not be perfect and may need some further adjustments, but those are not required to get full points in this question.

(a) (6 Points) Recreate the graphs (and layout) shown in Figure 3 using baseR. Include your R code and the resulting graphs. Hint: You can create a new line via `\n` without any extra spaces before/after `\n`.

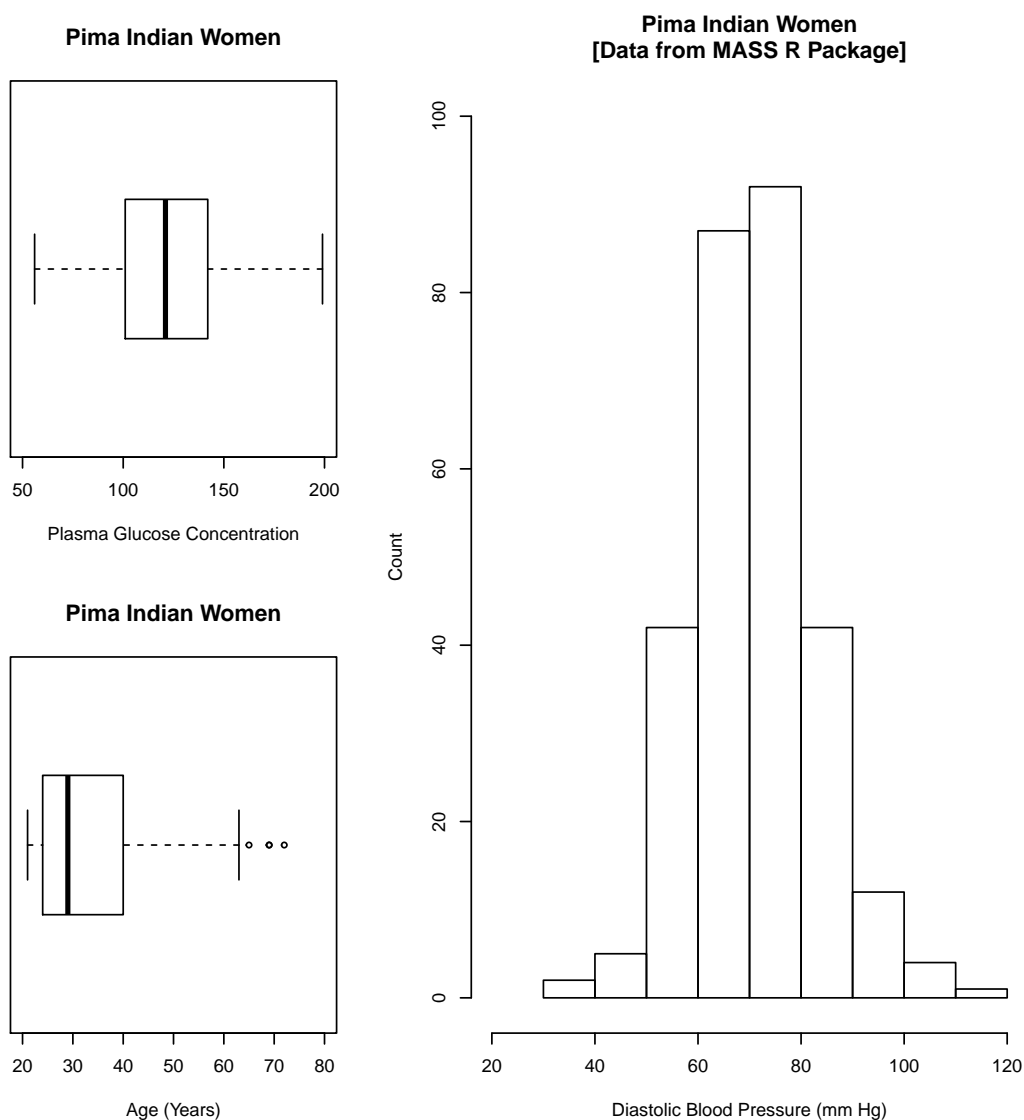


Figure 2: Graphs created with baseR.

- (b) (6 Points) Recreate the graph shown in Figure 3 using `ggplot2`. Include your R code and the resulting graph. Note: You have to change some of the labels and adjust the grid lines so they do not run through the middle of some of your bars. Find suitable help pages or information on [stackoverflow](https://stackoverflow.com).

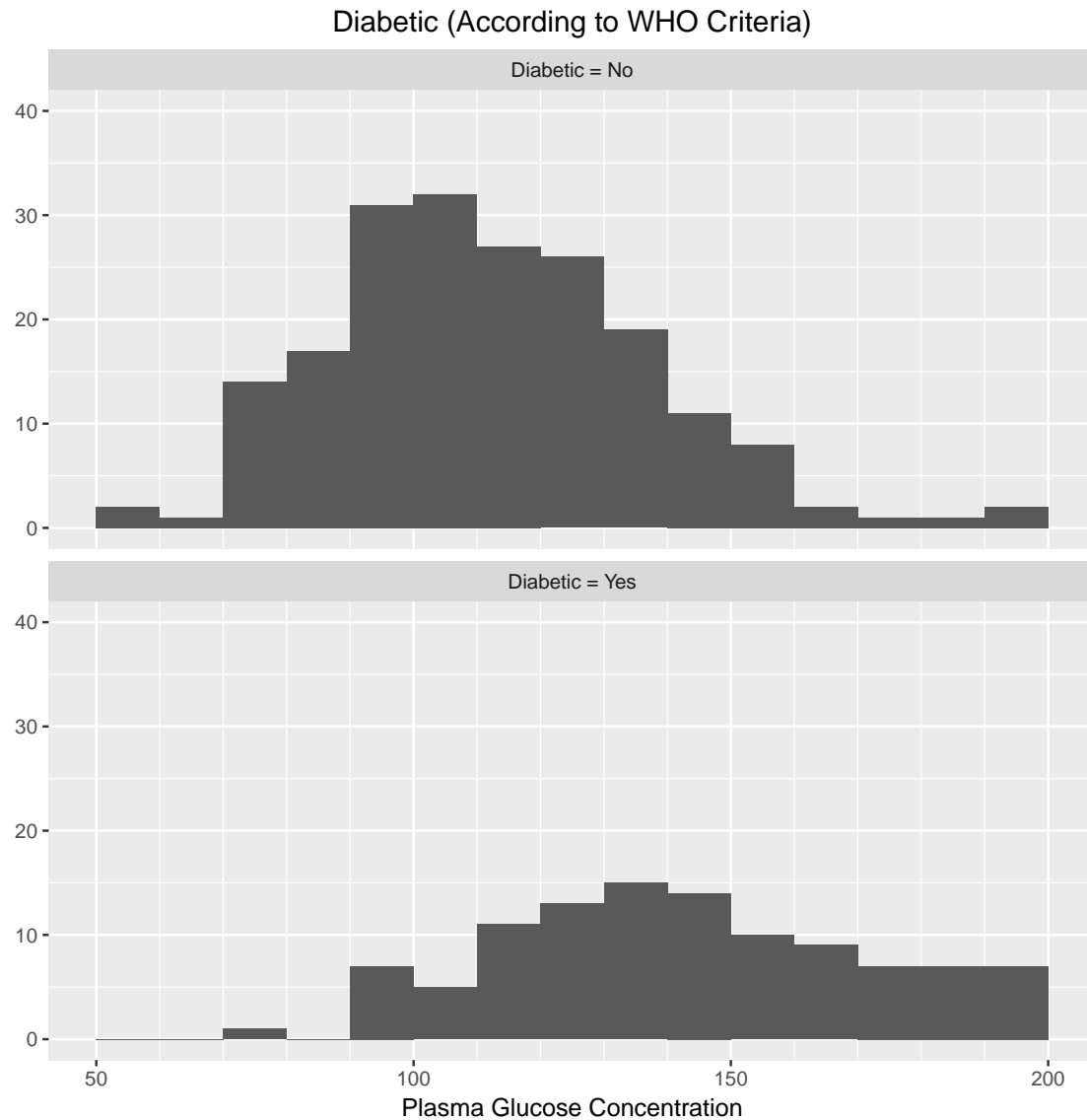


Figure 3: Graph created with *ggplot2*.

General Instructions

- (i) Create a single html or pdf document, using R Markdown, Sweave, or knitr. You only have to submit this one document.
- (ii) Include a title page that contains your name, your A-number, the number of the assignment, the submission date, and any other relevant information.
- (iii) Start your answers to each main question on a new page (continuing with the next part of a question on the same page is fine). Clearly label each question and question part.
- (iv) Before you submit your homework, check that you follow all recommendations from Google's R Style Guide (see <https://google.github.io/styleguide/Rguide.xml>). Moreover, make sure that your R code is consistent, i.e., that you use the same type of assignments and the same type of quotes throughout your entire homework.
- (v) Give credit to external sources, such as stackoverflow or help pages. Be specific and include the full URL where you found the help (or from which help page you got the information). Consider R code from such sources as “legacy code or third-party code” that does not have to be adjusted to Google's R Style (even though it would be nice, in particular if you only used a brief code segment).
- (vi) **Not following the general instructions outlined above will result in point deductions!**
- (vii) For general questions related to this homework, please use the corresponding discussion board in Canvas! I will try to reply as quickly as possible. Moreover, if one of you knows an answer, please post it. It is fine to refer to web pages and R commands, but do not provide the exact R command with all required arguments or which of the suggestions from a stackoverflow web page eventually worked for you! This will be the task for each individual student!
- (viii) Submit your single html or pdf file via Canvas by the submission deadline. Late submissions will result in point deductions as outlined on the syllabus.