

```

1 # --- -----
2 # --- Descriptive Statistics in R with the IRIS dataset
3 # ---
4 # --- V1 March 2021, D.Benninger - initial version for BINA-FS21
5 # --- V2 March 2021, dbf --- some typos corrected
6 # --- V3 January 2022, dbf --- some minor corrections (ggplot) for CAS BIA12
7 # ---
8 # --- Libraries: ggplot2, car
9 # ---
10 # --- Data: iris >> standard R dataset
11 # ---
12 # --- Links
13 # --- https://www.r-bloggers.com/2020/01/descriptive-statistics-in-r/
14 # ---
15 # ---
16 # --- -----
17
18 # PACKAGES installieren, falls nicht vorhanden
19 if(!"ggplot2" %in% rownames(installed.packages())) install.packages("ggplot2")
20 if(!"car" %in% rownames(installed.packages())) install.packages("car")
21
22 # Packages laden
23 library("ggplot2")
24 library("car")
25
26 # Set WORKING Directory
27 setwd(choose.dir())
28
29
30 # --- DATA OVERVIEW
31 # load the iris dataset and renamed it dat
32 dat <- iris
33
34 # --- PREVIEW of the dataset and its structure
35 # first observations
36 head(dat)
37
38 # --- Sepal.Length Sepal.Width Petal.Length Petal.Width Species
39 # --- 1 5.1 3.5 1.4 0.2 setosa
40 # --- 2 4.9 3.0 1.4 0.2 setosa
41 # --- 3 4.7 3.2 1.3 0.2 setosa
42 # --- 4 4.6 3.1 1.5 0.2 setosa
43 # --- 5 5.0 3.6 1.4 0.2 setosa
44 # --- 6 5.4 3.9 1.7 0.4 setosa
45
46 # structure of dataset
47 str(dat)
48
49 # --- 'data.frame': 150 obs. of 5 variables:
50 # --- $ Sepal.Length: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
51 # --- $ Sepal.Width : num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
52 # --- $ Petal.Length: num 1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
53 # --- $ Petal.Width : num 0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
54 # --- $ Species : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 1
55 1 1 ...
56
57 # --- MINIMUM and MAXIMUM
58 # -----
59 min(dat$Sepal.Length)
60 # --- [1] 4.3
61
62 max(dat$Sepal.Length)
63 # --- [1] 7.9
64
65 # alternative the range() function
66 rng <- range(dat$Sepal.Length)
67 rng
68 # --- [1] 4.3 7.9
69
70 rng[1] # MIN -- rng = name of the object specified above
71 # --- [1] 4.3
72 rng[2] # MAX

```

```

73 # --- [1] 7.9
74
75 # --- RANGE
76 max(dat$Sepal.Length) - min(dat$Sepal.Length)
77 # --- [1] 3.6
78
79
80
81 # --- MEAN, MEDIAN and QUANTILE
82 # -----
83 mean(dat$Sepal.Length)
84 # --- [1] 5.843333
85
86 median(dat$Sepal.Length)
87 # --- [1] 5.8
88
89 quantile(dat$Sepal.Length, 0.5)
90 # --- 50%
91 # --- 5.8
92
93 # --- 1st and 3rd QUARTILE
94 # -----
95 quantile(dat$Sepal.Length, 0.25) # first quartile
96 # --- 25%
97 # --- 5.1
98 quantile(dat$Sepal.Length, 0.75) # third quartile
99 # --- 75%
100 # --- 6.4
101
102 # --- other QUANTILES
103 quantile(dat$Sepal.Length, 0.4) # 4th decile
104 # --- 40%
105 # --- 5.6
106 quantile(dat$Sepal.Length, 0.98) # 98th percentile
107 # --- 98%
108 # --- 7.7
109
110 # --- Interquartile Range
111 IQR(dat$Sepal.Length)
112
113 # or alternatively with the quantile function
114 quantile(dat$Sepal.Length, 0.75) - quantile(dat$Sepal.Length, 0.25)
115 # --- 75%
116 # --- 1.3
117
118
119
120 # --- Standard DEVIATION and VARIANCE
121 # -----
122 sd(dat$Sepal.Length) # standard deviation
123 # --- [1] 0.8280661
124 var(dat$Sepal.Length) # variance
125 # --- [1] 0.6856935
126
127 # Tipp: compute the standard deviation (or variance) of multiple variables at the
128 # same time
129 # use lapply() with the appropriate statistics as second argument:
130 lapply(dat[, 1:4], sd)
131 # --- $Sepal.Length
132 # --- [1] 0.8280661
133 # ---
134 # --- $Sepal.Width
135 # --- [1] 0.4358663
136 # ---
137 # --- $Petal.Length
138 # --- [1] 1.765298
139 # ---
140 # --- $Petal.Width
141 # --- [1] 0.7622377
142
143 #
144 -----
145 ----

```

```

143 # --- SUMMARY
144 #
-----
145 summary(dat)
146 # --- Sepal.Length Sepal.Width Petal.Length Petal.Width
147 # --- Min. :4.300 Min. :2.000 Min. :1.000 Min. :0.100
148 # --- 1st Qu.:5.100 1st Qu.:2.800 1st Qu.:1.600 1st Qu.:0.300
149 # --- Median :5.800 Median :3.000 Median :4.350 Median :1.300
150 # --- Mean :5.843 Mean :3.057 Mean :3.758 Mean :1.199
151 # --- 3rd Qu.:6.400 3rd Qu.:3.300 3rd Qu.:5.100 3rd Qu.:1.800
152 # --- Max. :7.900 Max. :4.400 Max. :6.900 Max. :2.500
153 # --- Species
154 # --- setosa :50
155 # --- versicolor:50
156 # --- virginica :50
157
158
159 # Tipp: if you need these descriptive statistics by group use the by() function:
160 by(dat, dat$Species, summary)
161 # --- dat$Species: setosa
162 # --- Sepal.Length Sepal.Width Petal.Length Petal.Width
163 # --- Min. :4.300 Min. :2.300 Min. :1.000 Min. :0.100
164 # --- 1st Qu.:4.800 1st Qu.:3.200 1st Qu.:1.400 1st Qu.:0.200
165 # --- Median :5.000 Median :3.400 Median :1.500 Median :0.200
166 # --- Mean :5.006 Mean :3.428 Mean :1.462 Mean :0.246
167 # --- 3rd Qu.:5.200 3rd Qu.:3.675 3rd Qu.:1.575 3rd Qu.:0.300
168 # --- Max. :5.800 Max. :4.400 Max. :1.900 Max. :0.600
169 # --- Species
170 # --- setosa :50
171 # --- versicolor: 0
172 # --- virginica : 0
173 # ---
174 # ---
175 # ---
176 # --- -----
177 # --- dat$Species: versicolor
178 # --- Sepal.Length Sepal.Width Petal.Length Petal.Width Species
179 # --- Min. :4.900 Min. :2.000 Min. :3.00 Min. :1.000 setosa : 0
180 # --- 1st Qu.:5.600 1st Qu.:2.525 1st Qu.:4.00 1st Qu.:1.200 versicolor:50
181 # --- Median :5.900 Median :2.800 Median :4.35 Median :1.300 virginica : 0
182 # --- Mean :5.936 Mean :2.770 Mean :4.26 Mean :1.326
183 # --- 3rd Qu.:6.300 3rd Qu.:3.000 3rd Qu.:4.60 3rd Qu.:1.500
184 # --- Max. :7.000 Max. :3.400 Max. :5.10 Max. :1.800
185 # --- -----
186 # --- dat$Species: virginica
187 # --- Sepal.Length Sepal.Width Petal.Length Petal.Width
188 # --- Min. :4.900 Min. :2.200 Min. :4.500 Min. :1.400
189 # --- 1st Qu.:6.225 1st Qu.:2.800 1st Qu.:5.100 1st Qu.:1.800
190 # --- Median :6.500 Median :3.000 Median :5.550 Median :2.000
191 # --- Mean :6.588 Mean :2.974 Mean :5.552 Mean :2.026
192 # --- 3rd Qu.:6.900 3rd Qu.:3.175 3rd Qu.:5.875 3rd Qu.:2.300
193 # --- Max. :7.900 Max. :3.800 Max. :6.900 Max. :2.500
194 # --- Species
195 # --- setosa : 0
196 # --- versicolor: 0
197 # --- virginica :50
198 # ---
199 # ---
200 # ---
201
202
203
204 # --- HISTOGRAM
205 # -----
206 hist(dat$Sepal.Length)
207
208 # with ggplot2
209 ggplot(dat) +
210 aes(x = Sepal.Length) +
211 geom_histogram()
212
213 # by default, the number of bins is 30.

```

```

214 # you can change this value with geom_histogram(bins = 12) for instance.
215 ggplot(dat) +
216   aes(x = Sepal.Length) +
217   geom_histogram(bins = 12)
218
219 # --- BOXPLOT
220 # -----
221 boxplot(dat$Sepal.Length)
222
223 # Boxplots are even more informative when presented side-by-side
224 # for comparing and contrasting distributions from two or more groups.
225 # For instance, we compare the length of the sepal across the different species:
226 boxplot(dat$Sepal.Length ~ dat$Species)
227
228 # with ggplot2
229 ggplot(dat) +
230   aes(x = Species, y = Sepal.Length) +
231   geom_boxplot()
232
233 # --- SCATTERPLOT
234 # -----
235 plot(dat$Sepal.Length, dat$Petal.Length)
236
237 # with ggplot2
238 ggplot(dat) +
239   aes(x = Sepal.Length, y = Petal.Length) +
240   geom_point()
241
242 # and (colored) with differentiating the points according to a factor (i.e. species)
243 ggplot(dat) +
244   aes(x = Sepal.Length, y = Petal.Length, colour = Species) +
245   geom_point() +
246   scale_color_hue()
247
248 # --- DENSITY plot
249 # -----
250 plot(density(dat$Sepal.Length))
251
252 # -----
253 # --- Additional Graphical Representations
254 # -----
255 plot(iris)
256 ggplot(iris, aes(Petal.Length, Petal.Width, color = Species)) + geom_point()
257
258 # --- Attribute Statistical Key Values
259 # --- Explore Individual Variables
260 boxplot(iris$Sepal.Length, main="Sepal.Length")
261 boxplot(iris$Petal.Length, main="Petal.Length")
262
263 par(mfrow=c(1,2))
264 boxplot(iris$Petal.Length, iris$Petal.Width, main="Petal Measures")
265 boxplot(iris$Sepal.Length, iris$Sepal.Width, main="Sepal Measures")
266
267 par(mfrow=c(1,1))
268 boxplot(iris$Sepal.Length, iris$Sepal.Width, iris$Petal.Length, iris$Petal.Width,
269 main="IRIS Measures", ylab="length/width in cm")
270 boxplot(Sepal.Length ~ Species, data=iris, xlab="Species", ylab="Sepal.Length")
271
272 # --- Histograms and Density plots
273 hist(iris$Sepal.Length, breaks=10)
274 hist(iris$Sepal.Width, breaks=10)
275 hist(iris$Sepal.Length+iris$Sepal.Width)
276
277 plot(density(iris$Petal.Length))
278 plot(density(iris$Petal.Width))
279
280 # --- SCATTERPLOT
281 # --- Explore Multiple Variables

```

```

286 plot(iris$Sepal.Length,iris$Sepal.Width, xlim=c(0,10), ylim=c(0,10), main="Iris Data")
287
288 plot(iris$Sepal.Length, iris$Sepal.Width, col=iris$Species, pch=as.numeric(iris$
Species), main="Iris Data Scatterplot")
289
290 library(scatterplot3d)
291 scatterplot3d(iris$Petal.Width,iris$Sepal.Length,iris$Sepal.Width)
292 scatterplot3d(iris$Petal.Width,iris$Sepal.Length,iris$Sepal.Width, pch=21, bg=c("red"
,"green3","blue") [unclass(iris$Species)])
293
294 # --- Smooth SCATTERPLOT
295 smoothScatter(iris$Sepal.Length,iris$Sepal.Width)
296
297 # --- HEAT MAP
298 distM <- as.matrix(dist(iris[,1:4]), method="euclidean")
299 heatmap(distM)
300
301 # --- Other Visualization (cross reference, species)
302 pairs(iris)
303
304 quickplot(Sepal.Length, Sepal.Width, data=iris, facets= Species ~.,color = Species)
305

```