

# Übung 4

## Explorative Datenanalyse und Visualisierung

Wintersemester 2019

S. Döhler, B. Nedic (FBMN, h\_da)

---

**Name:** Roman Kessler

**Punkte:**

---

**Aufgabe 13.** Verwenden Sie den R-Datensatz `faithful` (dieser Datensatz ist im base package enthalten). Das Merkmal `eruptions` enthält die Länge von Eruptionen (in Minuten) des Old Faithful Geysirs im Yellowstone National Park in Wyoming, USA. Sie sollen die Verteilung der Eruptionen analysieren. Achten Sie bei allen Grafiken darauf, dass sowohl der Titel der Grafik als auch die Beschriftungen der Achsen möglichst aussagekräftig sind.

- (a) Laden Sie die Daten und erstellen Sie eine five-point summary.
- (b) Erstellen Sie einen Box-Plot und interpretieren Sie diesen.
- (c) Erstellen Sie einen Stem and leaf-Plot und interpretieren Sie diesen.
- (d) Plotten Sie die empirische Verteilungsfunktion.
- (e) Plotten Sie das Histogramm (inklusive eines "rug-plots") und interpretieren Sie dieses. Experimentieren Sie mit verschiedenen Klassenbreiten. Wie müssen Sie die Optionen einstellen, damit das Histogramm als Dichteschätzer geplottet wird? Wie erhält man ein ähnliches Bild wie beim stem and leaf-Plot?
- (f) Plotten Sie ein gleitendes Histogramm in das Histogramm mit ein. Experimentieren Sie mit verschiedenen Bandbreiten. Welche Bandbreite liefert ein "gutes" Ergebnis?
- (g) Wiederholen Sie Teil (f) mit dem Gauß- und Epanechnikov-Kernen.
- (h) Vergleichen Sie die verschiedenen Methoden. Was haben Sie über die Daten gelernt? Welche Methode(n) würden Sie einem Anwender empfehlen?

### Lösung

```
library(base)
df = data.frame(faithful)
```

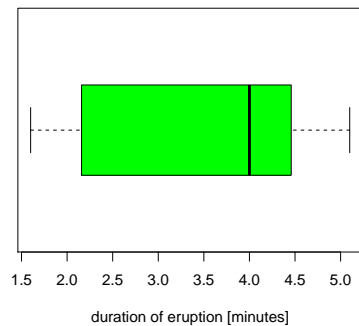
Zu (a):

Durch die Funktion "summary(dataFrame)" erhalten wir die 5 gesuchten Punkte und den arithmetischen Mittelwert:

```
summary(df)
#>      eruptions      waiting
#>  Min.      :1.600   Min.      :43.0
#> 1st Qu.:2.163   1st Qu.:58.0
#>  Median :4.000   Median :76.0
#>   Mean   :3.488   Mean   :70.9
#> 3rd Qu.:4.454   3rd Qu.:82.0
#>   Max.   :5.100   Max.   :96.0
```

Zu (b):

```
boxplot(df$eruptions, col = "green", xlab = "duration of eruption [minutes]", horizontal = T)
```



Wir sehen, dass die Dauer der Eruptionen zwischen etwa 1,5 und 5 Minuten ist, der Median jedoch bei etwa 4 Minuten liegt, was auf eine schiefe Verteilung der Daten schließen lässt.

Zu (c):

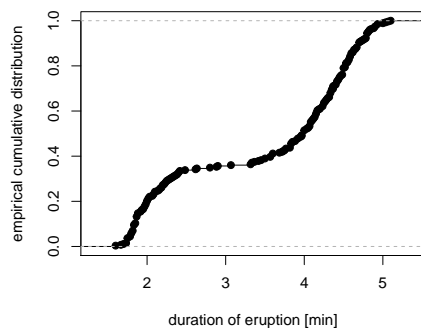
```
stem(df$eruptions)
#>
#> The decimal point is 1 digit(s) to the left of the |
#>
#> 16 | 070355555588
#> 18 | 00002223333333557777777888822335777888
#> 20 | 00002223378800035778
#> 22 | 0002335578023578
#> 24 | 00228
#> 26 | 23
#> 28 | 080
#> 30 | 7
#> 32 | 2337
#> 34 | 250077
#> 36 | 0000823577
#> 38 | 2333335582225577
```

```
#> 40 | 0000003357788888002233555577778
#> 42 | 03335555778800233333555577778
#> 44 | 02222335557780000000023333357778888
#> 46 | 00002333577000000023578
#> 48 | 00000022335800333
#> 50 | 0370
```

Der Stem-and-Leaf Plot gibt schon eine bessere Verteilung der Daten wieder. Wir sehen eine gewisse Dichotomie: Es gibt einige kurze Ausbrüche (um den Bereich von 2 Minuten) und viele längere Ausbrüche (um den Bereich von 3,5 bis 5 Minuten).

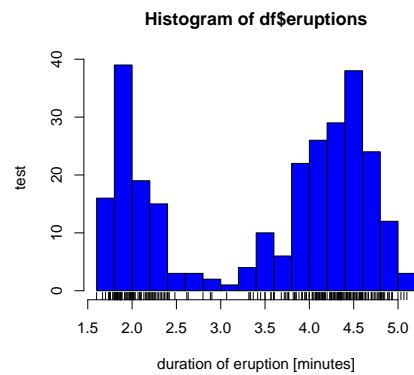
Zu (d):

```
plot(ecdf(df$eruptions), xlab = "duration of eruption [min]", ylab = "empirical cumulative c
```



Zu (e):

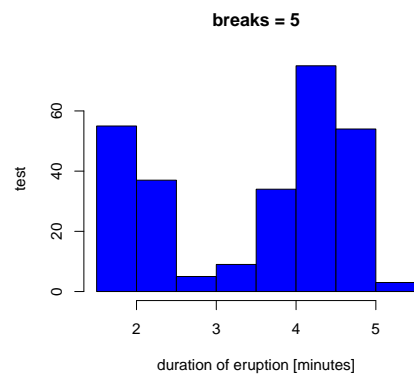
```
{
hist(df$eruptions, breaks = 20, col = "blue", xlab = "duration of eruption [minutes]", ylab = "frequency")
rug(df$eruptions)
}
```



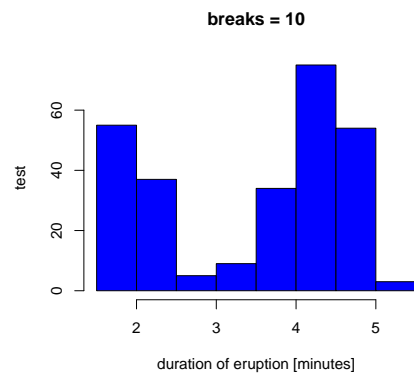
Das Histogramm (mit dem Rugplot) bestätigt noch einmal die Vermutung über die dichotome Verteilung der Daten.

Im folgenden experimentieren wir mit der Breite der “Bins”.

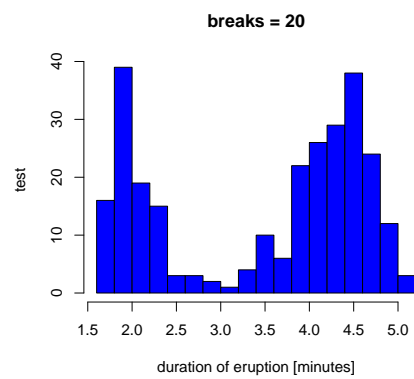
```
hist(df$eruptions, breaks = 5, col = "blue", xlab = "duration of eruption [minutes]", ylab =
```



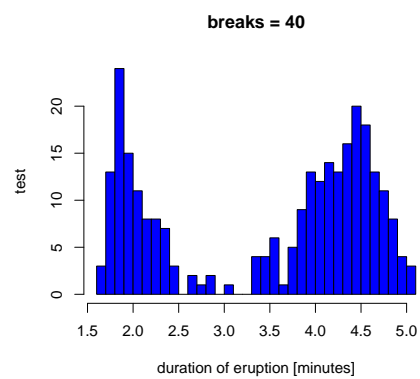
```
hist(df$eruptions, breaks = 10, col = "blue", xlab = "duration of eruption [minutes]", ylab =
```



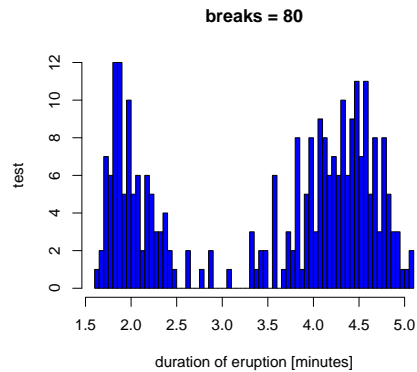
```
hist(df$eruptions, breaks = 20, col = "blue", xlab = "duration of eruption [minutes]", ylab
```



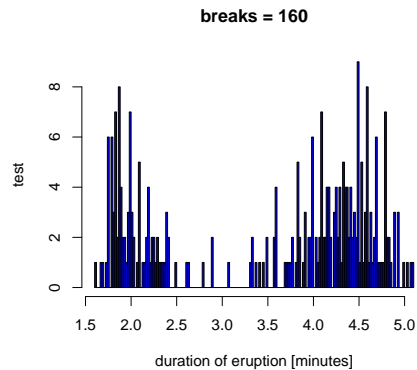
```
hist(df$eruptions, breaks = 40, col = "blue", xlab = "duration of eruption [minutes]", ylab
```



```
hist(df$eruptions, breaks = 80, col = "blue", xlab = "duration of eruption [minutes]", ylab = "test")
```



```
hist(df$eruptions, breaks = 160, col = "blue", xlab = "duration of eruption [minutes]", ylab = "test")
```



Wir sehen

1.: das Argument “breaks” (wenn wir ein Integer übergeben um die Anzahl der Zellen zu definieren), erst mal nur ein Vorschlag für die Histogramm Funktion darstellt. Bei einem Argument (z.B. breaks = 5), wird dieser Vorschlag ignoriert, da das Programm (laut Dokumentation: “pretty”) einen Wert festlegt, mit welchem ich als Betrachter die Verteilung der Daten besser dargestellt bekomme. Meine Eingabe (breaks = 5) scheint wohl eine ungünstiges Ergebnis zu liefern, welches mir eine andere Verteilung implizieren würde.

2.: Wähle ich die Bin-Anzahl zu hoch, liefert mir das keine “brauchbaren” Informationen über die generelle Verteilung der Daten. Eine Bin Größe, die der Gesamtzahl der Datenpunkte nahe kommt, ist somit weniger sinnvoll.

Zu (f):

```
library(aplpack)
slider.hist(df$eruptions, breaks = 5, col = "blue", xlab = "duration of eruption [minutes]").
#> Error in tkconfigure(sc, variable = slider1): konnte Funktion "tkconfigure" nicht finden
#?hist()
```

Zu (g):

Zu (h):

**Anmerkungen/Korrektur**

---