

Datenanalyse mit R

Schummelzettel

Kürzel für Schummelzettel

d – Ein Datensatz, in R ein Data Frame
v – ein einzelner Vektor, z.B. `d$Spaltenname`

Arbeitsverzeichnis

`getwd()`
 Zeigt das aktuelle Arbeitsverzeichnis an.
`setwd("C://dateipfad/arbeitsverzeichnis")`
 Verändert das aktuelle Arbeitsverzeichnis.

Datensatzverarbeitung

Einlesen xlsx-Paket (Excel-Dateien)
odbc-Paket (SQL-Datenbanken)

`d <- read.table("dateiname.txt")`
 Liest Textdatei und speichert diese in **d**.
`d <- read.csv("dateiname.csv")`
 Liest .csv-Datei mit , als Spaltentrenner und speichert diese in **d**.

! Mit dem Parameter **encoding** kann die Zeichencodierung verändert werden
 Bsp.: `encoding="UTF-8"` oder `encoding="latin1"`

`d <- read.csv2("dateiname.csv")`
 Liest .csv-Datei mit ; als Spaltentrenner und speichert diese in **d**.
`load(file = "dateiname.RData")`
 Liest .RData-Datei, lädt die darin gespeicherten Variablen. Überschreibt gleichlautende Variablen im Arbeitsplatz.
`d <- readRDS(file = "dateiname.rds")`
 Liest .rds-Datei, speichert das darin enthaltene R-Objekt in **d**.

Exportieren

`write.table(d, "dateiname.txt")`
 Schreibt Textdatei aus **d** und speichert diese in die Datei `dateiname.txt`.
`write.csv(d, file = "dateiname.csv")`
 Schreibt .csv-Datei aus **d** und speichert diese in die Datei `dateiname.csv`.

! Mit Parameter **row.names = FALSE** wird keine Spalte mit Durchnummerierung der Zeilen erzeugt

`save(d, file = "dateiname.RData")`
 Speichert **d** in die Datei `dateiname.Rdata` zum späteren Laden mit `load`.
`saveRDS(d, file = "dateiname.rds")`
 Speichert einzelnes R-Objekt **d** in die Datei `dateiname.rds` zum späteren Laden per `readRDS()`.

Datensatzuntersuchung/-vorbereitung

<code>nrow(d)</code>	Zeilenanzahl	<code>quantile(v, probs=x)</code>	Quantile von Spalte v ausgeben. Welche Quantile ausgegeben werden sollen wird mit dem Vektor x festgelegt.
<code>ncol(d)</code>	Spaltenanzahl	<code>factor(v)</code>	Spaltenformat zu Faktoren umwandeln. Damit werden Kategorien als solche gekennzeichnet.
<code>dim(d)</code>	Zeilenanzahl und Spaltenanzahl	<code>levels(v)</code>	Ausgabe der Kategorien bei Faktor-Format
<code>colnames(d)</code>	Spaltennamen	<code>length(unique(v))</code>	Anzahl der Kategorien
<code>head(d)</code>	erste 6 Zeilen	<code>table(v)</code>	Überblick wie viele Elemente pro Kategorie vorkommen
<code>tail(d)</code>	letzte 6 Zeilen	<code>plot(v)</code>	Diagramm für die zugehörige Spalte v . Diagrammart hängt von den Format der Spalte ab.
<code>View(d)</code>	Übersicht gesamter Datensatz	<code>qqnorm(v)</code> oder <code>EnvStats::qqPlot(v)</code>	Q-Q-Plot der Normalverteilung; <code>ggPlot()</code> stammt aus dem Paket <code>EnvStats</code>
		! Bei einem großen Datensatz ist ein Absturz von RStudio möglich.	
<code>str(d)</code>	Struktur des Datensatzes (Spaltennamen, Format der Elemente pro Spalte, erste Elemente jeder Spalte)		
<code>summary(d)</code>	Gibt einen zusammenfassenden Überblick für jede Spalte. Output unterscheidet sich je nach Typ der Spalte.		

Maße

	Lagemaße	Streuemaße
	Frage: Wo liegen die Werte?	Frage: Wie unterschiedlich sind die Werte?
Modus	<code>which.max(table(v))</code> <code>x <- v[complete.cases(v)]</code> <code>which.max(table(x))</code> Ohne Berücksichtigung von NAs	Varianz <code>var(v)</code> <code>var(v, na.rm = TRUE)</code> Ohne Berücksichtigung von NAs
Median	<code>median(v)</code> <code>median(v, na.rm = TRUE)</code> Ohne Berücksichtigung von NAs	Standardabweichung <code>sd(v)</code> <code>sd(v, na.rm = TRUE)</code> Ohne Berücksichtigung von NAs
Arithmetisches Mittel	<code>mean(v)</code> <code>mean(v, na.rm = TRUE)</code> Ohne Berücksichtigung von NAs	

Verteilungen

	Zufallsvariable	Dichtefunktion	Quantil
Normalverteilung	<code>rnorm</code>	<code>dnorm</code>	<code>qnorm</code>
Gleichverteilung	<code>runif</code>	<code>dunif</code>	<code>qunif</code>
Exponentialverteilung	<code>rexp</code>	<code>dexp</code>	<code>qexp</code>



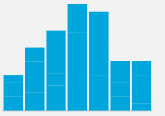
Histogramme

Ziel: Verteilung der Stichprobe grafisch darstellen

`hist(v)` Histogramm aus dem Vektor **v** erstellen

Sinnvolle Vorgänge:

- Anzahl der Balken verändern
`hist(v, breaks = zahl)`
- Datensatz in Subgruppen aufteilen
`d_anteil_eins <- subset(d, bedingung)`
`d_anteil_zwei <- subset(d, !bedingung)`



! Achten Sie darauf, dass an der Grenze der Bedingung keine Datenzeilen verloren gehen.

Histogramme über alle Variablen:

Hmisc - Paket

`hist(d)` Alle Variablen aus dem Datensatz **d** gleichzeitig ansprechen und pro Spalte ein Histogramm erstellen.

! Voraussetzung: alle Daten sind numerisch

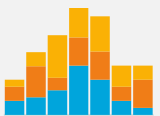
`Hmisc::hist.data.frame(d)` Alle Variablen aus dem Datensatz **d** gleichzeitig ansprechen und pro Spalte ein Histogramm erstellen.

Segmentierungsanalyse mit Histogrammen

ggplot2 - Paket
viridis - Paket

Additives Histogramm

`ggplot2::ggplot(d, aes(x=v1, fill = v2)) + geom_histogram(binwidth = zahl_breite)`



Additives Histogramm erstellen mit Datensatz **d** mit den Inhalten aus Spalte **v1**, wobei sich die Balken farblich nach Spalte **v2** aufeinander legen.

Separiertes Histogramm

`ggplot2::ggplot(d, aes(x=v1, fill = v2)) + geom_histogram(binwidth = zahl_breite, position = "dodge")`



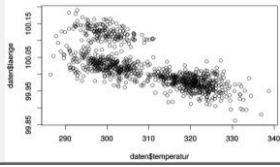
Separiertes Histogramm erstellen mit Datensatz **d** mit den Inhalten aus Spalte **v1**, wobei sich die Balken farblich nach Spalte **v2** separieren.

Streudiagramme und Matrixplot

Ziel: Variablen Abhängigkeiten bei numerischen Daten entdecken.

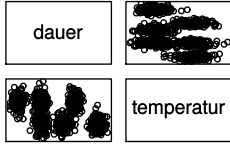
Streudiagramm (Scatter Plot / XY-Plot)

`plot(x = v1, y = v2)`
Stellt die Elemente der Spalte v1 mit den Elementen der Spalte v2 in Abhängigkeit dar.



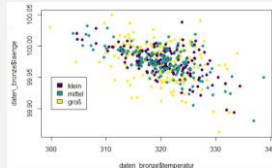
Matrixplot

`spalten <- c(...)`
`plot(d[, spalten])`
Stellt jede Kombination von Variablen getrennt dar. Hilfreich ist eine Vorauswahl von Spalten.



Segmentierung mit Streudiagramm

`plot(x = v1, y = v2, pch = 19)`
`d1 <- subset(d, bedingung1)`
`d2 <- subset(d, bedingung2)`
...



`lines(d1$v1, d1$v2, type = "p", pch=19)`
`lines(d2$v1, d2$v2, type = "p", pch=19)`
...
Stellt Streudiagramm mit den Elementen aus Spalte v1 und Spalte v2 in Abhängigkeit dar. Der Datensatz wird nach Bedingung aufgeteilt und jeder Datensatzanteil bekommt eine andere Farbe.

Parameter `type = "p"` sorgt dafür, dass Punkte und keine Linien erzeugt werden.
Parameter `pch = 19` sorgt dafür, dass die Punkte ausgefüllt sind.

Korrelationsanalyse

Ziel: Objektive Kennzahl für die Abhängigkeit von Variablen ermitteln.

`cor(v1, v2)`

Korrelationskoeffizient zwischen Spalte v1 und Spalte v2.

! Setzen sie außerdem den Parameter `use` auf `use = "pairwise.complete.obs"`. Damit werden fehlende Werte in einer Spalte nicht berücksichtigt.

`spalten <- c(...)`
`cor(d[, spalten])`

Korrelationskoeffizienten -matrix von allen ausgewählten Spalten erstellen.

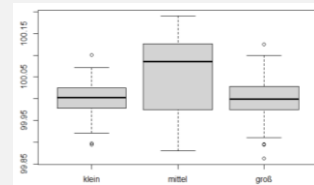
Grafische Aufbereitung:
`corrplot`-Paket oder `ggplot2`-Paket mit `Ggally`-Paket

Boxplots und Multi-Vari Chart

Ziel: Variablen Abhängigkeiten bei kategorialen und stetigen Daten entdecken.

Boxplot (Kastendiagramm)

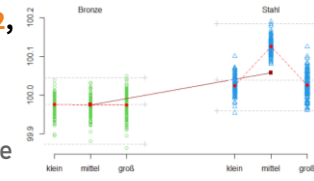
`plot(x = v1, y = v2)`
Stellt die Elemente der Spalte v1 mit den Elementen der Spalte v2 in Abhängigkeit dar. x muss dabei eine Faktorvariable erhalten.



Multi-Vari Chart

qualityTools - Paket

`qualityTools::mvPlot(v1, v2, v3, FUN = ...)`
Zwei kategoriale Daten aus den Spalten v1 und v2 in Abhängigkeit zu einer Spalte v3 mit stetigen Variablen.



Hypothesentest

Ziel: Beobachtungen von Eigenschaften der Grundgesamtheit statistisch verifizieren.

Vorgehensweise

Formulierung von Nullhypothese H_0 und Alternativhypothese H_1

H_0	H_1	Einordnung
$\theta \geq \theta_0$	$\theta < \theta_0$	Einseitiger Test (linksseitig)
$\theta \leq \theta_0$	$\theta > \theta_0$	Einseitiger Test (rechtsseitig)
$\theta = \theta_0$	$\theta \neq \theta_0$ ($\theta < \theta_0$ oder $\theta > \theta_0$)	Zweiseitiger Test

1.

2. Auswahl eines geeigneten statistischen Tests

3. Festlegen der Irrtumswahrscheinlichkeit (meist 5%, 1%, oder 0.1%)

4. Durchführung des Tests: Vergleich von berechneter Testgröße mit dem kritischen Wert aus der Verteilung der Teststatistik

5. Treffen der Testentscheidung:
→ Falls berechneter Wert den kritischen Wert über/unterschreitet ist H_0 abzulehnen
→ Andernfalls ist H_0 nicht abzulehnen und H_1 ist plausibel

Fehler 1. und 2. Art

	H_0 ist richtig	H_0 ist falsch
H_0 wird nicht verworfen	Richtige Entscheidung Wahrscheinlichkeit: $1 - \alpha$	Fehler 2. Art Wahrscheinlichkeit: β
H_0 wird verworfen	Fehler 1. Art Wahrscheinlichkeit: α	Richtige Entscheidung Wahrscheinlichkeit: $1 - \beta$

ANOVA

Ziel: Mehrere Kategorien miteinander vergleichen im Hinblick der Mittelwerte.

`a <- aov(y ~ k, data = d)`
`summary(a)`
a - Ergebnis der ANOVA
y - Zu prüfende Variable
k - Kategorivariable
d - Datensatz

Interpretation:
`Pr(>F)` klein -> y ist über k hinweg im Mittel nicht gleich

Lineare Regression

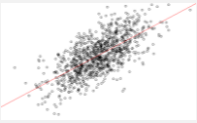
Ziel: Lineare Beziehung zwischen Variablen bestimmen.

Einfache lineare Regression:

Optimale Parameter a und b schätzen:
 $Y = a + bX + \epsilon$

`lm(y ~ x, data = d)`

y / x - Die Variable y aus Datensatz d soll durch die Variable x aus Datensatz d erklärt werden.



Interpretation:

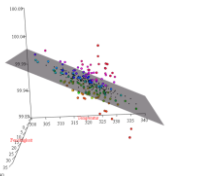
`coefficients[1]` bzw. **Intercept** → a
`coefficients[2]` → b

Multiple lineare Regression:

Optimale Parameter a, b_1, b_2, \dots schätzen:
 $Y = a + b_1X_1 + b_2X_2 + \dots + \epsilon$

`lm(y ~ x1 + x2 + ..., data = d)`

y / x1 / x2 - Die Variable y aus Datensatz d soll durch die Variablen x1 und x2 (...) aus Datensatz d erklärt werden.



Interpretation:

`coefficients[1]` bzw. **Intercept** → a
`coefficients[2]` → b_1
`coefficients[3]` → b_2
...

Statistische Verifikation

`summary(lm(...))`

Signifikanz der Regression anzeigen.

Bereich	Interpretation
Coefficients: ...	Beschreibung der Koeffizienten
	Estimate - Spalte gibt die geschätzten Parameter an Pr(> t) - Spalte berechnete p-Werte von Hypothesentest „Parameter = 0“ * - Sterne Geben Signifikanzniveau an
Residuals: ...	Einschätzung wie gut die Regression die Daten beschreibt
	Es gilt: umso < Residuen, desto besser beschreibt die Regressionsgerade die Datenpunkte
Multiple R-squared: ...	%-der Varianz der Daten, die durch die Einflussvariablen beschrieben werden kann = 1 → Daten perfekt beschrieben
F-statistic: ...	Hypothesentest: alle Einflussparameter zusammen 0

Konfidenzintervalle

Ziel: Bereich bestimmen, in dem der tatsächliche Wert mit hoher Wahrscheinlichkeit basierend auf der Stichprobe unter Modellannahmen liegt

Konfidenzniveau	$1 - \alpha$
Irrtumswahrscheinlichkeit	α Gibt an, wie viel Prozent aller Fälle der wahre Wert nicht im Konfidenzintervall liegt

