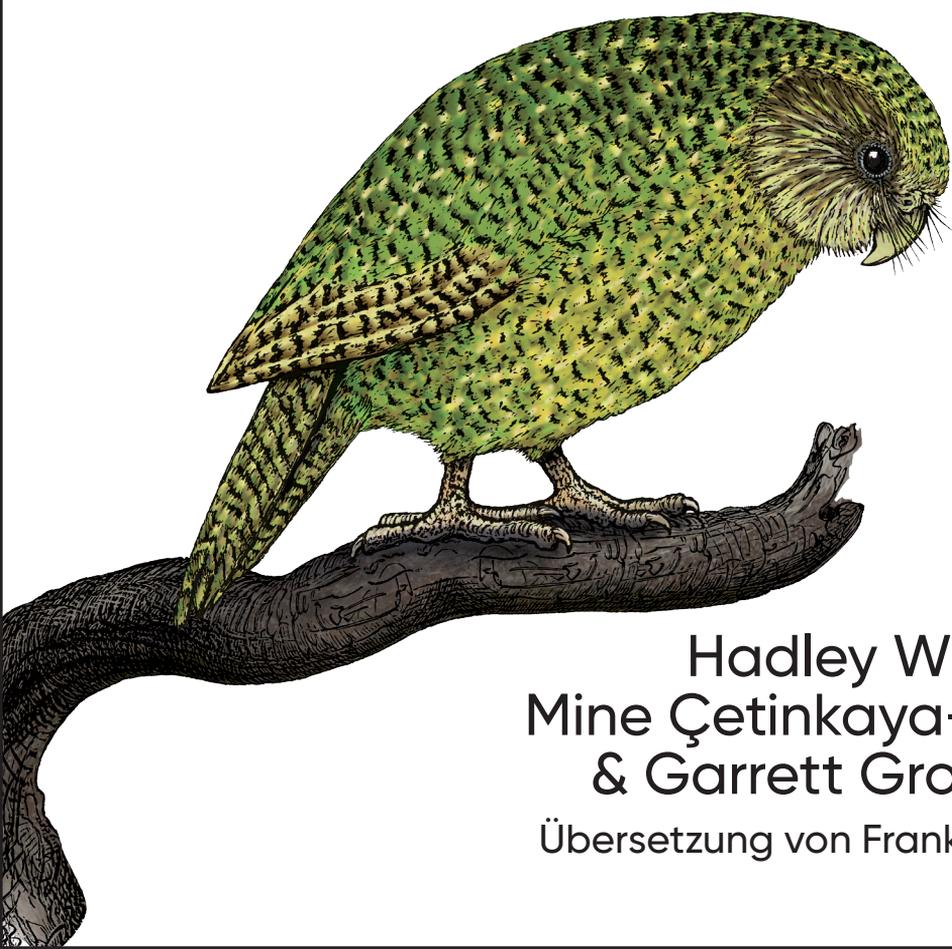


O'REILLY®

2. Auflage
aktualisiert
und erweitert

R für Data Science

Daten importieren, bereinigen, umformen
und visualisieren



Hadley Wickham,
Mine Çetinkaya-Rundel
& Garrett Golemund

Übersetzung von Frank Langenau

Daten importieren

Einführung

Mit Daten zu arbeiten, die von den R-Paketen bereitgestellt werden, ist eine großartige Möglichkeit, die Tools der Data Science kennenzulernen. Doch an einem gewissen Punkt möchten Sie das Gelernte auch auf Ihre eigenen Daten anwenden. Dieses Kapitel beschäftigt sich nun mit den Grundlagen, Datendateien in R einzulesen.

Speziell konzentriert es sich darauf, wie man einfache, rechteckige Textdateien einliest. Los geht es mit einem praktischen Hinweis für den Umgang mit Features wie Spaltennamen, Typen und fehlenden Daten. Dann erfahren Sie, wie man mehrere Dateien auf einmal einliest und Daten aus R in eine Datei schreibt. Schließlich lernen Sie, wie Sie Dataframes in R manuell zusammenbauen.

Voraussetzungen

In diesem Kapitel lernen Sie, lineare Dateien mit dem Paket `readr` in R zu laden. Dieses Paket ist Teil des Kern-Tidyverse:

```
library(tidyverse)
```

Daten aus einer Datei lesen

Zu Beginn konzentrieren wir uns auf den gebräuchlichsten Typ einer rechteckigen Datendatei: CSV, was als Abkürzung für *Comma-Separated Values* (kommagetrennte Werte) steht. Das folgende Beispiel zeigt, wie eine einfache CSV-Datei aussieht. Die erste Zeile, häufig auch *Header-Zeile* oder Überschriftenzeile genannt, gibt die Spaltennamen an, und in den folgenden sechs Zeilen sind die Daten enthalten. Die Spalten werden durch Kommata voneinander getrennt.

```

Student ID,Full Name,favourite.food,mealPlan,AGE
1,Sunil Huffmann,Strawberry yoghurt,Lunch only,4
2,Barclay Lynn,French fries,Lunch only,5
3,Jayendra Lyne,N/A,Breakfast and lunch,7
4,Leon Rossini,Anchovies,Lunch only,
5,Chidiegwu Dunkel,Pizza,Breakfast and lunch,five
6,Güvenç Attila,Ice cream,Lunch only,6

```

Tabelle 7-1 stellt dieselben Daten als Tabelle dar.

Tabelle 7-1: Daten aus der Datei `students.csv` in Form einer Tabelle

Student ID	Full Name	favourite.food	mealPlan	AGE
1	Sunil Huffmann	Strawberry yoghurt	Lunch only	4
2	Barclay Lynn	French fries	Lunch only	5
3	Jayendra Lyne	N/A	Breakfast and lunch	7
4	Leon Rossini	Anchovies	Lunch only	NA
5	Chidiegwu Dunkel	Pizza	Breakfast and lunch	five
6	Güvenç Attila	Ice cream	Lunch only	6

Diese Datei können wir mit `read_csv()` in R einlesen. Das erste Argument ist das wichtigste: der Pfad zur Datei. Man kann sich den Pfad als die Adresse der Datei vorstellen: Die Datei heißt `students.csv` und »wohnt« im Ordner `data`.

```

students <- read_csv("data/students.csv")
#> Rows: 6 Columns: 5
#> — Column specification —————
#> Delimiter: ","
#> chr (4): Full Name, favourite.food, mealPlan, AGE
#> dbl (1): Student ID
#>
#> i Use `spec()` to retrieve the full column specification for this data.
#> i Specify the column types or set `show_col_types = FALSE` to quiet this
message.

```

Der obige Code funktioniert, wenn Sie die Datei `students.csv` in einem Ordner `data` in Ihrem Projekt abgelegt haben. Die Datei `students.csv` können Sie von <https://oreil.ly/GDubbb> herunterladen oder direkt von der angegebenen URL lesen:

```
students <- read_csv("https://pos.it/r4ds-students-csv")
```

Die Funktion `read_csv()` gibt eine Meldung zurück mit der Anzahl der Datenzeilen und -spalten, dem verwendeten Trennzeichen und den Spaltenspezifikationen (Namen der Spalten, organisiert nach dem Typ der in den Spalten enthaltenen Daten). Außerdem erfahren Sie, wie sich die vollständige Spaltenspezifikation abrufen lässt und wie Sie diese Meldung unterdrücken können. Die Meldung ist integraler Bestandteil von `readr`, und wir kommen im Abschnitt »Spaltentypen steuern« auf Seite 132 darauf zurück.

Praktischer Ratschlag

Nachdem Sie Daten eingelesen haben, besteht der erste Schritt in der Regel darin, sie in bestimmter Weise umzuwandeln, um sie für die weitere Analyse leichter verarbeiten zu können. Sehen wir uns in diesem Sinne noch einmal die students-Daten an:

```
students
#> # A tibble: 6 × 5
#>   `Student ID` `Full Name` favourite.food mealPlan AGE
#>   <dbl> <chr> <chr> <chr> <chr>
#> 1 1 Sunil Huffmann Strawberry yoghurt Lunch only 4
#> 2 2 Barclay Lynn French fries Lunch only 5
#> 3 3 Jayendra Lyne N/A Breakfast and lunch 7
#> 4 4 Leon Rossini Anchovies Lunch only <NA>
#> 5 5 Chidiegwu Dunkel Pizza Breakfast and lunch five
#> 6 6 Güvenç Attila Ice cream Lunch only 6
```

In der Spalte `favourite.food` stehen verschiedene Lebensmittel sowie die Zeichenfolge »N/A«, die ein richtiges NA sein sollte, das R als *not available* (nicht verfügbar) erkennt. Das ist etwas, das wir mit dem Argument `na` angehen können. Standardmäßig erkennt die Funktion `read_csv()` nur leere Zeichenfolgen ("") in diesem Datenset als NA-Werte. Wir möchten aber, dass sie auch den String "N/A" erkennt:

```
students <- read_csv("data/students.csv", na = c("N/A", ""))

students
#> # A tibble: 6 × 5
#>   `Student ID` `Full Name` favourite.food mealPlan AGE
#>   <dbl> <chr> <chr> <chr> <chr>
#> 1 1 Sunil Huffmann Strawberry yoghurt Lunch only 4
#> 2 2 Barclay Lynn French fries Lunch only 5
#> 3 3 Jayendra Lyne <NA> Breakfast and lunch 7
#> 4 4 Leon Rossini Anchovies Lunch only <NA>
#> 5 5 Chidiegwu Dunkel Pizza Breakfast and lunch five
#> 6 6 Güvenç Attila Ice cream Lunch only 6
```

Sicherlich haben Sie bemerkt, dass die Spalten `Student ID` und `Full Name` von Backticks umgeben sind. Das liegt daran, dass sie Leerzeichen enthalten und damit die üblichen Regeln von R für Variablenamen verletzen. Um auf diese Variablen zu verweisen, müssen Sie sie in Backticks (` `) einschließen:

```
students |>
  rename(
    student_id = `Student ID`,
    full_name = `Full Name`
  )
#> # A tibble: 6 × 5
#>   student_id full_name favourite.food mealPlan AGE
#>   <dbl> <chr> <chr> <chr> <chr>
#> 1 1 Sunil Huffmann Strawberry yoghurt Lunch only 4
#> 2 2 Barclay Lynn French fries Lunch only 5
```

```

#> 3      3 Jayendra Lyne   <NA>           Breakfast and lunch 7
#> 4      4 Leon Rossini   Anchovies       Lunch only           <NA>
#> 5      5 Chidiegwu Dunkel Pizza           Breakfast and lunch five
#> 6      6 Güvenç Attila  Ice cream       Lunch only           6

```

Als alternativer Ansatz bietet sich die Funktion `janitor::clean_names()` an, um mithilfe einer Heuristik alle Namen auf einmal in Snake Case umzuwandeln:¹

```

students |> janitor::clean_names()
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <chr>          <chr>
#> 1     1 Sunil Huffmann Strawberry yoghurt Lunch only      4
#> 2     2 Barclay Lynn   French fries   Lunch only      5
#> 3     3 Jayendra Lyne <NA>          Breakfast and lunch 7
#> 4     4 Leon Rossini   Anchovies       Lunch only      <NA>
#> 5     5 Chidiegwu Dunkel Pizza           Breakfast and lunch five
#> 6     6 Güvenç Attila  Ice cream       Lunch only      6

```

Nach dem Einlesen der Daten ist es häufig ebenfalls erforderlich, die Variablentypen zu betrachten. Zum Beispiel ist `meal_plan` eine kategoriale Variable mit einem bekannten Satz möglicher Werte, die in R als Faktor dargestellt werden sollte.

```

students |>
  janitor::clean_names() |>
  mutate(meal_plan = factor(meal_plan))
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <fct>          <chr>
#> 1     1 Sunil Huffmann Strawberry yoghurt Lunch only      4
#> 2     2 Barclay Lynn   French fries   Lunch only      5
#> 3     3 Jayendra Lyne <NA>          Breakfast and lunch 7
#> 4     4 Leon Rossini   Anchovies       Lunch only      <NA>
#> 5     5 Chidiegwu Dunkel Pizza           Breakfast and lunch five
#> 6     6 Güvenç Attila  Ice cream       Lunch only      6

```

Beachten Sie, dass die Werte in der Variablen `meal_plan` gleich geblieben sind, aber der Variablentyp, der unter dem Variablennamen angegeben ist, hat sich von Zeichen (`<chr>`) in Faktor (`<fct>`) geändert. Kapitel 16 geht näher auf Faktoren ein.

Bevor Sie diese Daten analysieren, werden Sie wahrscheinlich die Spalten `age` und `id` bereinigen wollen. Derzeit ist `age` eine Zeichenvariable, weil eine der Beobachtungen als Zahlwort `five` ausgeschrieben ist statt als Ziffer 5. Wie sich derartige Probleme korrigieren lassen, besprechen wir ausführlich in Kapitel 20.

```

students <- students |>
  janitor::clean_names() |>
  mutate(
    meal_plan = factor(meal_plan),
    age = parse_number(if_else(age == "five", "5", age))
  )

```

¹ Das Paket `janitor` (<https://oreil.ly/J8GX>) ist nicht im Tidyverse enthalten, bietet aber praktische Funktionen für die Datenbereinigung und funktioniert auch gut in Daten-Pipelines, die `|>` verwenden.

```

students
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <fct>          <dbl>
#> 1     1 Sunil Huffmann  Strawberry yoghurt Lunch only      4
#> 2     2 Barclay Lynn   French fries    Lunch only      5
#> 3     3 Jayendra Lyne <NA>           Breakfast and lunch 7
#> 4     4 Leon Rossini  Anchovies      Lunch only      NA
#> 5     5 Chidiegwu Dunkel Pizza          Breakfast and lunch 5
#> 6     6 Güvenç Attila  Ice cream      Lunch only      6

```

Neu ist hier die Funktion `if_else()`, die drei Argumente hat. Das erste Argument `test` sollte ein logischer Vektor sein. Das Ergebnis enthält den Wert des zweiten Arguments, `yes`, wenn `test` gleich `TRUE` ist, und den Wert des dritten Arguments, `no`, wenn der Test `FALSE` liefert. Hier sagen wir: Wenn `age` den String "five" enthält, mache "5" daraus, und wenn nicht, bleibt `age`, wie es war. Mehr über `if_else()` und logische Vektoren lernen Sie in Kapitel 12.

Andere Argumente

Es gibt noch eine Reihe weiterer wichtiger Argumente, die wir erwähnen müssen, und sie lassen sich besser vorführen, wenn wir uns zunächst einen praktischen Trick ansehen: Die Funktion `read_csv()` kann Textzeichenfolgen lesen, die Sie erzeugt und wie eine CSV-Datei formatiert haben:

```

read_csv(
  "a,b,c
  1,2,3
  4,5,6"
)
#> # A tibble: 2 × 3
#>   a     b     c
#>   <dbl> <dbl> <dbl>
#> 1     1     2     3
#> 2     4     5     6

```

Normalerweise verwendet `read_csv()` die erste Zeile der Daten für die Spaltennamen, was eine gängige Konvention ist. Es ist auch nicht ungewöhnlich, dass am Anfang der Datei einige Zeilen mit Metadaten enthalten sind. Die ersten `n` Zeilen können Sie mit `skip = n` überspringen, oder Sie können mit `comment = "#"` alle Zeilen auslassen, die beispielsweise mit `#` beginnen:

```

read_csv(
  "The first line of metadata
  The second line of metadata
  x,y,z
  1,2,3", skip = 2
)
#> # A tibble: 1 × 3
#>   x     y     z
#>   <dbl> <dbl> <dbl>
#> 1     1     2     3

```

```

read_csv(
  "# A comment I want to skip
  x,y,z
  1,2,3", comment = "#"
)
#> # A tibble: 1 × 3
#>       x     y     z
#>   <dbl> <dbl> <dbl>
#> 1     1     2     3

```

Es kann auch sein, dass die Daten keine Spaltennamen haben. Mit `col_names = FALSE` weisen Sie `read_csv()` an, die erste Zeile nicht als Überschriftenzeile (Header) zu verarbeiten, sondern sie stattdessen sequenziell von X1 bis Xn zu beschriften:

```

read_csv(
  "1,2,3
  4,5,6",
  col_names = FALSE
)
#> # A tibble: 2 × 3
#>       X1    X2    X3
#>   <dbl> <dbl> <dbl>
#> 1     1     2     3
#> 2     4     5     6

```

Alternativ können Sie `col_names` einen Zeichenvektor übergeben, der für die Spaltennamen verwendet wird:

```

read_csv(
  "1,2,3
  4,5,6",
  col_names = c("x", "y", "z")
)
#> # A tibble: 2 × 3
#>       x     y     z
#>   <dbl> <dbl> <dbl>
#> 1     1     2     3
#> 2     4     5     6

```

Wenn Sie diese Argumente kennen, sind Sie in der Lage, die meisten in der Praxis vorkommenden CSV-Dateien einzulesen. (Für die übrigen Varianten müssen Sie Ihre `.csv`-Datei sorgfältig inspizieren und die Dokumentation für die vielen anderen Argumente von `read_csv()` studieren.)

Andere Dateitypen

Sobald Sie `read_csv()` beherrschen, ist es sehr einfach, die anderen Funktionen von `readr` zu verwenden. Sie müssen lediglich wissen, welche Funktion jeweils infrage kommt:

`read_csv2()`

Liest Dateien ein, deren Felder durch Semikola (;) statt durch Kommata (,) getrennt sind. Derartige Dateien sind in Ländern üblich, in denen das Komma als Dezimaltrennzeichen dient.

`read_tsv()`

Liest Dateien ein, deren Felder durch Tabulatoren getrennt sind.

`read_delim()`

Liest Dateien mit einem beliebigen Trennzeichen ein, wobei versucht wird, das Trennzeichen automatisch zu erraten, wenn Sie es nicht angeben.

`read_fwf()`

Liest Dateien mit Feldern fester Breite ein. Die Felder können Sie mit `fwf_widths()` durch ihre Breite oder mit `fwf_positions()` durch ihre Positionen spezifizieren.

`read_table()`

Liest eine gebräuchliche Variation von Dateien mit fester Breite ein, wobei die Spalten durch Leerzeichen getrennt sind.

`read_log()`

Liest Protokolldateien im Apache-Stil ein.

Übungen

1. Mit welcher Funktion würden Sie eine Datei einlesen, deren Felder durch | getrennt sind?
2. Welche anderen Argumente außer `file`, `skip` und `comment` haben die Funktionen `read_csv()` und `read_tsv()` gemein?
3. Was sind die wichtigsten Argumente der Funktion `read_fwf()`?
4. Manchmal enthalten CSV-Dateien Zeichenfolgen mit Kommata. Um Probleme zu vermeiden, müssen diese Kommata in Anführungszeichen eingeschlossen werden, und zwar in einfache (') oder doppelte ("). Standardmäßig geht `read_csv()` davon aus, dass es sich um doppelte Anführungszeichen (") handelt. Welches Argument müssen Sie bei `read_csv()` angeben, um den folgenden Text in einen Dataframe einzulesen?

```
"x,y\n1, 'a,b'"
```

5. Ermitteln Sie, was bei den folgenden Inline-CSV-Dateien jeweils nicht stimmt. Was passiert, wenn Sie den Code ausführen?

```
read_csv("a,b\n1,2,3\n4,5,6")  
read_csv("a,b,c\n1,2\n1,2,3,4")  
read_csv("a,b\n\"1")  
read_csv("a,b\n1,2\na,b")  
read_csv("a;b\n1;3")
```

6. Üben Sie, sich auf nicht syntaktische Namen im folgenden Dataframe zu beziehen, indem Sie
 - a. die Variable namens 1 extrahieren,
 - b. ein Streudiagramm von 1 gegen 2 erstellen,
 - c. eine neue Spalte namens 3 erzeugen, die 2 geteilt durch 1 ist,
 - d. die Spalten in one, two und three umbenennen:

```
annoying <- tibble(
  `1` = 1:10,
  `2` = `1` * 2 + rnorm(length(`1`))
)
```

Spaltentypen steuern

Da eine CSV-Datei keine Informationen über den Typ jeder Variablen enthält (d.h., ob sie einen logischen Wert, eine Zahl, eine Zeichenfolge usw. darstellt), versucht `readr`, den Typ zu erraten. Dieser Abschnitt beschreibt, wie das Erraten funktioniert, wie man einige häufige Probleme löst, die zum Scheitern führen, und wie man bei Bedarf die Spaltentypen selbst bereitstellen kann. Zum Schluss stellen wir noch einige allgemeine Strategien vor, die nützlich sind, wenn `readr` katastrophal versagt und Sie mehr Einblick in die Struktur Ihrer Datei benötigen.

Typen erraten

Das Paket `readr` verwendet eine Heuristik, um die Spaltentypen herauszufinden. Für jede Spalte holt es die Werte von 1.000 Zeilen² in gleichmäßigen Abständen von der ersten bis zur letzten Zeile und ignoriert dabei fehlende Werte. Anschließend arbeitet es die folgenden Fragen ab:

- Enthält sie nur F, T, FALSE oder TRUE (ohne Beachtung der Groß-/Kleinschreibung)? Wenn ja, handelt es sich um einen logischen Wert.
- Enthält sie nur Zahlen (z. B. 1, -4.5, 5e6, Inf)? Wenn ja, handelt es sich um eine Zahl.
- Entspricht sie dem Standard ISO8601? Wenn ja, handelt es sich um ein Datum oder um ein Datum mit Uhrzeit. (Auf Datums-/Zeitwerte kommen wir ausführlich in Kapitel 17 im Abschnitt »Datums-/Zeitwerte erzeugen« auf Seite 326 zurück.)
- Andernfalls muss es sich um eine Zeichenfolge handeln.

Dieses Verhalten können Sie mit diesem einfachen Beispiel nachvollziehen:

```
read_csv("
  logical,numeric,date,string
```

² Den Standardwert 1.000 können Sie mit dem Argument `guess_max` überschreiben.

```

TRUE,1,2021-01-15,abc
false,4.5,2021-02-15,def
T,Inf,2021-02-16,ghi
")
#> # A tibble: 3 × 4
#>   logical numeric date      string
#>   <lg1>    <dbl> <date>   <chr>
#> 1 TRUE      1  2021-01-15 abc
#> 2 FALSE    4.5 2021-02-15 def
#> 3 TRUE     Inf  2021-02-16 ghi

```

Diese Heuristik funktioniert gut, wenn das Datenset sauber ist, doch in der Praxis werden Sie auf eine Reihe von seltsamen und ungewöhnlichen Fehlern stoßen.

Fehlende Werte, Spaltentypen und Probleme

Eine Spaltenerkennung scheitert vor allem dann, wenn eine Spalte unerwartete Werte enthält. Dann bekommen Sie eine Zeichenspalte anstelle eines spezifischeren Typs. Eine der häufigsten Ursachen dafür ist ein fehlender Wert, der mit etwas anderem als dem von `readr` erwarteten `NA` erfasst wurde.

Nehmen Sie diese einfache einspaltige CSV-Datei als Beispiel:

```

simple_csv <- "
x
10
.
20
30"

```

Wenn wir die Datei ohne zusätzliche Argumente einlesen, wird `x` zu einer Zeichenspalte:

```

read_csv(simple_csv)
#> # A tibble: 4 × 1
#>   x
#>   <chr>
#> 1 10
#> 2 .
#> 3 20
#> 4 30

```

In diesem kleinen Datenset können Sie den fehlenden Wert – durch einen Punkt (`.`) dargestellt – leicht erkennen. Doch wie sieht es aus bei Tausenden von Zeilen mit nur wenigen fehlenden Werten, die durch Punkte dargestellt werden? Man könnte `readr` mitteilen, dass `x` eine numerische Spalte ist, und dann sehen, wo das Einlesen versagt. Hierfür weisen Sie dem Argument `col_types` eine benannte Liste zu, in der die Namen den Spaltennamen in der CSV-Datei entsprechen:

```

df <- read_csv(
  simple_csv,
  col_types = list(x = col_double())
)

```

```
#> Warning: One or more parsing issues, call `problems()` on your data frame for
#> details, e.g.:
#> dat <- vroom(...)
#> problems(dat)
```

Jetzt meldet `read_csv()`, dass es ein Problem gibt, und sagt uns, dass wir mit `problems()` mehr herausfinden können:

```
problems(df)
#> # A tibble: 1 × 5
#>   row col expected actual file
#>   <int> <int> <chr>   <chr> <chr>
#> 1     3     1 a double .   /private/tmp/RtmpAY1Sop/file392d445cf269
```

Wir erfahren nun, dass es ein Problem in Zeile 3, Spalte 1 gibt, wo `readr` einen Wert vom Typ `double` erwartet, aber einen Punkt (`.`) vorgefunden hat. Das legt nahe, dass dieses Datenset fehlende Werte mit einem Punkt kennzeichnet. Also setzen wir `na = "."`. Die automatische Typherleitung ist nun erfolgreich und liefert uns die numerische Spalte, die wir haben wollten:

```
read_csv(simple_csv, na = ".")
#> # A tibble: 4 × 1
#>       x
#>   <dbl>
#> 1    10
#> 2    NA
#> 3    20
#> 4    30
```

Spaltentypen

Im Paket `readr` können Sie aus insgesamt neun Spaltentypen wählen:

- `col_logical()` und `col_double()` lesen logische Werte und Realzahlen. Man benötigt sie nur selten (außer wie oben gezeigt), da `readr` normalerweise solche Typen automatisch erkennt.
- `col_integer()` liest Ganzzahlen. In diesem Buch unterscheiden wir nur selten zwischen Ganzzahlen und Gleitkommazahlen (`double`), da sie funktional äquivalent sind. Allerdings kann das explizite Lesen von Ganzzahlen gelegentlich nützlich sein, da sie gegenüber Gleitkommazahlen nur die Hälfte des Speichers belegen.
- `col_character()` liest Zeichenfolgen. Dies kann nützlich sein, wenn zum Beispiel eine Spalte einen numerischen Bezeichner verkörpert, d.h. eine lange Folge von Ziffern, die ein Objekt identifizieren, aber in mathematischen Operationen nicht sinnvoll sind. Beispiele hierfür sind Telefonnummern, Sozialversicherungsnummern, Kreditkartennummern usw.
- `col_factor()`, `col_date()` und `col_datetime()` erzeugen Faktoren, Datumswerte und Datums-/Zeitwerte. Mehr dazu erfahren Sie, wenn wir in den Kapiteln 16 und 17 auf diese Datentypen zu sprechen kommen.

- `col_number()` ist ein toleranter numerischer Parser, der nicht numerische Komponenten ignoriert und besonders nützlich ist für Währungen. Mehr dazu lesen Sie in Kapitel 13.
- `col_skip()` überspringt eine Spalte, die auch nicht in das Ergebnis aufgenommen wird. Dies kann zum Beispiel Zeit sparen, wenn Sie sehr große CSV-Dateien einlesen müssen, aber nur einige der Spalten benötigen.

Es ist auch möglich, die Standardspalte zu überschreiben, indem Sie von `list()` zu `cols()` wechseln und `.default` angeben:

```
another_csv <- "
x,y,z
1,2,3"

read_csv(
  another_csv,
  col_types = cols(.default = col_character())
)
#> # A tibble: 1 × 3
#>   x     y     z
#>   <chr> <chr> <chr>
#> 1 1     2     3
```

Eine andere nützliche Hilfsfunktion ist `cols_only()`, die nur die angegebene(n) Spalte(n) einliest:

```
read_csv(
  another_csv,
  col_types = cols_only(x = col_character())
)
#> # A tibble: 1 × 1
#>   x
#>   <chr>
#> 1 1
```

Daten aus mehreren Dateien einlesen

Manchmal sind Ihre Daten auf mehrere Dateien verteilt, anstatt in einer einzigen Datei enthalten zu sein. Ein Beispiel hierfür sind monatliche Umsatzdaten, wobei die Daten für jeden Monat in einer eigenen Datei liegen: `01-sales.csv` für Januar, `02-sales.csv` für Februar und `03-sales.csv` für März. Mit `read_csv()` können Sie diese Daten auf einmal einlesen und sie in einem einzelnen Dataframe übereinanderstapeln.

```
sales_files <- c("data/01-sales.csv", "data/02-sales.csv", "data/03-sales.csv")
read_csv(sales_files, id = "file")
#> # A tibble: 19 × 6
#>   file                month   year brand  item    n
#>   <chr>              <chr> <dbl> <dbl> <dbl> <dbl>
#> 1 data/01-sales.csv January 2019     1 1234     3
#> 2 data/01-sales.csv January 2019     1 8721     9
```

```
#> 3 data/01-sales.csv January 2019 1 1822 2
#> 4 data/01-sales.csv January 2019 2 3333 1
#> 5 data/01-sales.csv January 2019 2 2156 9
#> 6 data/01-sales.csv January 2019 2 3987 6
#> # ... with 13 more rows
```

Auch hier funktioniert der obige Code nur dann, wenn Sie die CSV-Dateien in einem Ordner *data* in Ihrem Projekt abgelegt haben. Diese drei Dateien können Sie von <https://oreil.ly/jVd8o>, <https://oreil.ly/RYsgM> und <https://oreil.ly/4uZOm> herunterladen oder sie wie folgt direkt einlesen:

```
sales_files <- c(
  "https://pos.it/r4ds-01-sales",
  "https://pos.it/r4ds-02-sales",
  "https://pos.it/r4ds-03-sales"
)
read_csv(sales_files, id = "file")
```

Das Argument *id* fügt eine neue Spalte namens *file* in den resultierenden Dataframe ein. In ihr ist die Datei angegeben, aus der die Daten stammen. Das ist vor allem dann hilfreich, wenn die einzulesenden Dateien keine Identifizierungsspalte haben, die es erlauben würde, die Beobachtungen zu ihren ursprünglichen Quellen zurückzuverfolgen.

Wenn Sie viele Dateien einlesen wollen, kann es recht umständlich sein, ihre Namen als Liste zu schreiben. Stattdessen können Sie die Basisfunktion `list.files()` verwenden, die Ihnen die gewünschten Dateien anhand eines Musters in den Dateinamen zusammensucht. Mehr über derartige Muster lernen Sie in Kapitel 15.

```
sales_files <- list.files("data", pattern = "sales\\.csv$", full.names = TRUE)
sales_files
#> [1] "data/01-sales.csv" "data/02-sales.csv" "data/03-sales.csv"
```

In eine Datei schreiben

Das Paket `readr` bringt ebenfalls zwei nützliche Funktionen mit, um Daten auf einen Datenträger zu schreiben: `write_csv()` und `write_tsv()`. Die wichtigsten Argumente dieser Funktionen sind *x* (der zu speichernde Dataframe) und *file* (der Ort, an dem die Datei zu speichern ist). Außerdem können Sie mit *na* festlegen, wie fehlende Werte gespeichert werden sollen, und mit *append*, ob Sie den Dataframe an eine vorhandene Datei anfügen wollen.

```
write_csv(students, "students.csv")
```

Lesen wir nun diese CSV-Datei erneut ein. Beachten Sie, dass die Informationen über die Variablentypen, die Sie eben eingerichtet haben, verloren gehen, wenn Sie den Dataframe als CSV speichern. Somit stehen Sie wieder am Anfang und müssen eine reine Textdatei einlesen:

```

students
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <fct>          <dbl>
#> 1     1 Sunil Huffmann  Strawberry yoghurt Lunch only      4
#> 2     2 Barclay Lynn   French fries    Lunch only      5
#> 3     3 Jayendra Lyne <NA>           Breakfast and lunch 7
#> 4     4 Leon Rossini   Anchovies      Lunch only      NA
#> 5     5 Chidiegwu Dunkel Pizza          Breakfast and lunch 5
#> 6     6 Güvenç Attila  Ice cream      Lunch only      6
write_csv(students, "students-2.csv")
read_csv("students-2.csv")
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <chr>          <dbl>
#> 1     1 Sunil Huffmann  Strawberry yoghurt Lunch only      4
#> 2     2 Barclay Lynn   French fries    Lunch only      5
#> 3     3 Jayendra Lyne <NA>           Breakfast and lunch 7
#> 4     4 Leon Rossini   Anchovies      Lunch only      NA
#> 5     5 Chidiegwu Dunkel Pizza          Breakfast and lunch 5
#> 6     6 Güvenç Attila  Ice cream      Lunch only      6

```

Deshalb sind CSV-Dateien nicht so recht geeignet, um Zwischenergebnisse zu speichern – Sie müssen die Spaltenspezifikation jedes Mal erneut durchführen, wenn Sie den Dataframe laden. Hierzu gibt es zwei Alternativen:

- `write_rds()` und `read_rds()` sind einheitliche Wrapper um die Basisfunktionen `readRDS()` und `saveRDS()`. Diese Funktionen speichern die Daten in dem R-eigenen Binärformat namens RDS. Wenn Sie also das Objekt zurückladen, dann laden Sie *genau das gleiche* R-Objekt zurück, das Sie gespeichert haben.

```

write_rds(students, "students.rds")
read_rds("students.rds")
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <fct>          <dbl>
#> 1     1 Sunil Huffmann  Strawberry yoghurt Lunch only      4
#> 2     2 Barclay Lynn   French fries    Lunch only      5
#> 3     3 Jayendra Lyne <NA>           Breakfast and lunch 7
#> 4     4 Leon Rossini   Anchovies      Lunch only      NA
#> 5     5 Chidiegwu Dunkel Pizza          Breakfast and lunch 5
#> 6     6 Güvenç Attila  Ice cream      Lunch only      6

```

- Das Paket `arrow` erlaubt es, Dateien im Datendateiformat Parquet zu lesen und zu schreiben. Dieses schnelle, binäre Dateiformat lässt sich über Programmiersprachen hinweg einsetzen. Auf das Paket `arrow` kommen wir ausführlich in Kapitel 22 zurück.

```

library(arrow)
write_parquet(students, "students.parquet")
read_parquet("students.parquet")
#> # A tibble: 6 × 5
#>   student_id full_name      favourite_food meal_plan      age
#>   <dbl> <chr>          <chr>          <fct>          <dbl>

```

```

#> 1      1 Sunil Huffmann  Strawberry yoghurt Lunch only      4
#> 2      2 Barclay Lynn   French fries      Lunch only      5
#> 3      3 Jayendra Lyne  NA                Breakfast and lunch 7
#> 4      4 Leon Rossini   Anchovies         Lunch only      NA
#> 5      5 Chidiegwu Dunkel Pizza              Breakfast and lunch 5
#> 6      6 Güvencü Attila Ice cream           Lunch only      6

```

Parquet ist in der Regel viel schneller als RDS und lässt sich auch außerhalb von R einsetzen. Allerdings erfordert es das Paket `arrow`.

Dateneingabe

Manchmal müssen Sie ein Tibble »manuell« zusammenstellen, indem Sie ein wenig Dateingabe in Ihrem R-Skript praktizieren. Hierbei helfen Ihnen zwei nützliche Funktionen, die sich darin unterscheiden, ob das Tibble spalten- oder zeilenorientiert ist. Die Funktion `tibble()` arbeitet spaltenorientiert:

```

tibble(
  x = c(1, 2, 5),
  y = c("h", "m", "g"),
  z = c(0.08, 0.83, 0.60)
)
#> # A tibble: 3 × 3
#>       x y       z
#>   <dbl> <chr> <dbl>
#> 1     1 h     0.08
#> 2     2 m     0.83
#> 3     5 g     0.6

```

Wenn die Daten spaltenorientiert angeordnet sind, lässt sich schwerer erkennen, welche Beziehungen zwischen den Zeilen bestehen. Eine Alternative ist also `tribble()`, kurz für *transposed tibble*, mit dem Sie Ihre Daten zeilenweise anordnen können. Die Funktion `tribble()` ist auf die Dateneingabe im Code angepasst: Spaltenüberschriften beginnen mit einer Tilde (~), und die Einträge werden durch Kommata getrennt. Dadurch ist es möglich, kleine Datenmengen in einer gut lesbaren Form anzuordnen:

```

tribble(
  ~x, ~y, ~z,
  1, "h", 0.08,
  2, "m", 0.83,
  5, "g", 0.60
)
#> # A tibble: 3 × 3
#>       x       y       z
#>   <chr> <dbl> <dbl>
#> 1     1     h     0.08
#> 2     2     m     0.83
#> 3     5     g     0.6

```

Zusammenfassung

In diesem Kapitel haben Sie gelernt, wie Sie CSV-Dateien mit `read_csv()` laden und Ihre eigene Dateneingabe mit `tibble()` und `tribble()` realisieren. Es wurde gezeigt, wie CSV-Dateien funktionieren, auf welche Probleme Sie möglicherweise stoßen und wie Sie diese lösen können. In diesem Buch werden wir noch mehrmals auf den Datenimport zurückkommen: In Kapitel 20 laden Sie Daten aus Excel und Google Sheets, in Kapitel 21 aus Datenbanken, in Kapitel 22 aus Parquet-Dateien, in Kapitel 23 von JSON und in Kapitel 24 von Websites.

Wir sind fast am Ende dieses Abschnitts des Buchs angelangt, aber es gibt noch ein wichtiges letztes Thema zu behandeln: wie man Hilfe bekommt. Im nächsten Kapitel erfahren Sie, wo Sie am besten nach Hilfe suchen können und wie Sie ein Reprex erstellen, um die Chancen auf gute Hilfe zu maximieren, und Sie erhalten einige allgemeine Ratschläge dazu, wie Sie in der Welt von R auf dem Laufenden bleiben können.

Einführung	13
-------------------------	-----------

Teil I Gesamtbild

1 Datenvisualisierung	29
Einführung	29
Erste Schritte	30
Zielsetzung	32
ggplot2-Aufrufe	42
Verteilungen visualisieren	42
Beziehungen visualisieren	46
Diagramme speichern	53
Häufige Probleme	54
Zusammenfassung	55
2 Workflow: Grundlagen	57
Grundlagen der Codierung	57
Kommentare	58
Was macht einen Namen aus?	59
Funktionen aufrufen	60
Zusammenfassung	62
3 Datentransformation	63
Einführung	63
Voraussetzungen	63
Zeilen	66
Spalten	71
Die Pipe	76

Gruppen	78
Fallstudie: Aggregate und Stichprobengröße	85
Zusammenfassung	87
4 Workflow: Programmierstil	89
Namen	90
Leerzeichen	90
Pipes	91
ggplot2	93
Abschnittskommentare	93
Übungen	94
Zusammenfassung	94
5 Datenaufbereitung	95
Einführung	95
Aufbereitete Daten	96
Daten länger machen	99
Daten breiter machen	108
Zusammenfassung	112
6 Workflow: Skripte und Projekte	113
Skripte	113
Projekte	117
Übungen	122
Zusammenfassung	122
7 Daten importieren	125
Einführung	125
Daten aus einer Datei lesen	125
Spaltentypen steuern	132
Daten aus mehreren Dateien einlesen	135
In eine Datei schreiben	136
Dateneingabe	138
Zusammenfassung	139
8 Workflow: Hilfe abrufen	141
Google ist Ihr Freund	141
Ein Reprex erstellen	142
Sich selbst einbringen	144
Zusammenfassung	144

Teil II Visualisieren

9	Datenvisualisierung	147
	Einführung	147
	Ästhetische Zuordnungen	148
	Geometrische Objekte	152
	Facetten	158
	Statistische Transformationen	160
	Positionsanpassungen	164
	Koordinatensysteme	168
	Die geschichtete Grammatik der grafischen Darstellung	170
	Zusammenfassung	171
10	Explorative Datenanalyse	173
	Einführung	173
	Fragen	174
	Variation	175
	Ungewöhnliche Werte	179
	Kovariation	182
	Muster und Modelle	192
	Zusammenfassung	195
11	Kommunikation	197
	Einführung	197
	Beschriftungen	198
	Anmerkungen	200
	Skalen	205
	Themen	220
	Layout	223
	Zusammenfassung	227

Teil III Transformieren

12	Logische Vektoren	231
	Einführung	231
	Vergleiche	232
	Boolesche Algebra	236
	Zusammenfassungen	239
	Bedingte Transformationen	242
	Zusammenfassung	246

13 Zahlen	247
Einführung	247
Zahlen erzeugen	247
Zähler	248
Numerische Transformationen	250
Zahlen in Bereiche aufteilen	256
Allgemeine Transformationen	257
Numerische Zusammenfassungen	261
Zusammenfassung	267
14 Strings	269
Einführung	269
Einen String erzeugen	270
Viele Strings aus Daten erstellen	272
Daten aus Strings extrahieren	275
Buchstaben	281
Nicht englischer Text	283
Zusammenfassung	286
15 Reguläre Ausdrücke	287
Einführung	287
Muster – Grundlagen	288
Wichtige Funktionen	290
Details zu Mustern	295
Mustersteuerung	302
Praxis	304
Reguläre Ausdrücke an anderen Stellen	309
Zusammenfassung	311
16 Faktoren	313
Einführung	313
Faktoren-Basics	313
General Social Survey	315
Faktorreihenfolge ändern	316
Faktorlevels ändern	321
Geordnete Faktoren	323
Zusammenfassung	324
17 Datum und Uhrzeit	325
Einführung	325
Datums-/Zeitwerte erzeugen	326
Datums-/Zeitkomponenten	332

Zeiträume	339
Zeitzone n	343
Zusammenfassung	345
18 Fehlende Werte	347
Einführung	347
Explizit fehlende Werte	347
Implizit fehlende Werte	349
Faktoren und leere Gruppen	352
Zusammenfassung	355
19 Verknüpfungen	357
Einführung	357
Schlüssel	358
Grundlegende Verknüpfungen	363
Wie funktionieren Verknüpfungen?	369
Nicht-Gleichheitsverknüpfungen	375
Zusammenfassung	381

Teil IV Importieren

20 Tabellenkalkulationen	385
Einführung	385
Excel	385
Google Sheets	398
Zusammenfassung	403
21 Datenbanken	405
Einführung	405
Datenbankgrundlagen	406
Mit einer Datenbank verbinden	406
Grundlagen von dbplyr	409
SQL	411
Übersetzung von Funktionen	420
Zusammenfassung	423
22 Arrow	425
Einführung	425
Die Daten erhalten	426
Ein Datenset öffnen	426
Das Parquet-Format	428

Das Paket dplyr mit Arrow verwenden	430
Das Paket dbplyr mit Arrow verwenden	432
Zusammenfassung	433
23 Hierarchische Daten	435
Einführung	435
Listen	436
Verschachtelung beseitigen	440
Fallstudien	444
JSON	452
Zusammenfassung	456
24 Web-Scraping	457
Einführung	457
Ethische und rechtliche Anmerkungen zum Scraping	458
HTML-Grundlagen	460
Daten extrahieren	461
Die richtigen Selektoren finden	466
Alles zusammen	467
Dynamische Sites	472
Zusammenfassung	472

Teil V Programmieren

25 Funktionen	475
Einführung	475
Vektorfunktionen	476
Dataframe-Funktionen	482
Diagrammfunktionen	489
Stil	495
Zusammenfassung	496
26 Iteration	499
Einführung	499
Mehrere Spalten ändern	500
Mehrere Dateien lesen	509
Mehrere Ausgaben speichern	518
Zusammenfassung	523

27 Ein Praxisleitfaden zu Basis-R	525
Einführung	525
Mehrere Elemente mit [auswählen	526
Ein einzelnes Element mit \$ und [[auswählen	530
Familie der apply-Funktionen	533
for-Schleifen	535
Diagramme	537
Zusammenfassung	538

Teil VI Kommunizieren

28 Quarto	541
Einführung	541
Quarto-Grundlagen	542
Visueller Editor	545
Quelltexteditor	547
Codeblöcke	549
Bilder	553
Tabellen	557
Zwischenspeichern (Caching)	558
Fehlersuche	560
YAML-Header	561
Workflow	564
Zusammenfassung	566
29 Quarto-Formate	567
Einführung	567
Ausgabeoptionen	567
Dokumente	568
Präsentationen	569
Interaktivität	569
Websites und Bücher	572
Andere Formate	573
Zusammenfassung	573
Index	575

R für Data Science

Erfahren Sie, wie Sie mit R aus Ihren Daten Erkenntnisse und Einsichten gewinnen. Dieses Buch führt Sie in R und RStudio ein sowie in Tidyverse, eine Sammlung von R-Paketen, mit denen Data-Science-Aufgaben effektiv und zeitsparend erledigt werden können. Auch wenn Sie keine Programmiererfahrung haben, können Sie mit diesem aktualisierten Standardwerk schnell in die Praxis der Data Science einsteigen.

Sie lernen, Daten zu importieren, aufzubereiten, zu visualisieren und die Ergebnisse zu präsentieren. Darüber hinaus bekommen Sie einen umfassenden Überblick über den Data-Science-Zyklus und die Tools, die für die Detailarbeit erforderlich sind. Die zweite Auflage behandelt die neuesten Funktionen und Best Practices von Tidyverse und zeigt Ihnen in neu hinzugekommenen Kapiteln, wie Sie Daten aus Spreadsheets, Datenbanken und Websites nutzen. Zahlreiche Übungen unterstützen Sie dabei, das Gelernte praktisch auszuprobieren.

Themen des Buchs sind:

- **Visualisieren** – Erstellen Sie Diagramme für die Datenauswertung und die Kommunikation von Ergebnissen
- **Transformieren** – Erkunden Sie Variablentypen und die Werkzeuge, um mit ihnen zu arbeiten
- **Importieren** – Lesen Sie Daten in R ein und bringen Sie sie in eine für die Analyse geeignete Form
- **Programmieren** – Lernen Sie leistungsfähige R-Tools kennen, mit denen Sie Datenprobleme leichter lösen können
- **Kommunizieren** – Verwenden Sie Quarto, um Text, Code und Ergebnisse kombiniert darzustellen

»Dies ist eine außerordentlich gute Neuauflage des weltweit führenden Handbuchs für Data Science mit R. Alle, die mit Daten arbeiten, sollten es lesen.«

– Emma Rand
University of York, UK

Hadley Wickham ist Chief Scientist bei Posit und Mitglied der R Foundation. Er erstellt sowohl technische als auch kognitive Tools, die Data Science leichter, schneller und unterhaltsamer machen.

Mine Çetinkaya-Rundel ist Professor of the Practice und Director of Undergraduate Studies am Department of Statistical Science an der Duke University, North Carolina. Außerdem ist sie Developer Educator bei Posit.

Garrett Grolemond ist Statistiker, Lehrer und Director of Learning bei der Posit Academy, und er ist Autor von *Hands-On Programming with R* (O'Reilly).



www.dpunkt.de

Euro 54,90 (D)
ISBN 978-3-96009-234-6

plus⁺

Interesse am E-Book?
www.dpunkt.plus



Gedruckt in Deutschland
Papier aus nachhaltiger Waldwirtschaft
Mineralölfreie Druckfarben