

Constructing data frames (Gegevensframes construeren)

Constructing a data frame

In R-programmering is een "gegevensframe" een tweedimensionale structuur die gegevens opslaat in een tabelformaat van cel-rijen en cel-kolommen. In tegenstelling tot matrixstructuren hoeven de waarden die in een gegevensframe zijn opgeslagen niet allemaal van hetzelfde gegevenstype te zijn - ze kunnen waarden van elk gegevenstype bevatten. Dit betekent dat dataframes bijzonder veelzijdig zijn en de meest bruikbare datastructuur in R-programmering zijn.

Een dataframe kan in R Script worden gemaakt door een unieke identificatienaam naar keuze in de Code Editor te schrijven en vervolgens waarden toe te wijzen met behulp van de ingebouwde **data.frame()**-functie. Deze functie vereist vectoren met waarden als argumenten. Elke vector moet dezelfde lengte hebben, anders worden de waarden in extra cellen gerecycleerd om overeen te komen met de lengte van de langste vector.



Vergeet niet om het puntteken op te nemen in de functienamen **dataframe()** en **is.data.frame()**.

De waarden in elke vector verschijnen in afzonderlijke kolommen van het gegevensframe en elke kolom krijgt standaard de naam van de overeenkomstige vector als kolomnaam. Rijnamen worden standaard in oplopende volgorde vanaf één genummerd. Net als bij matrices kan de rangschikking van rijen en kolommen van dataframes worden getransponeerd met behulp van de functie **t()** - om rijen naar kolommen te schakelen.

Als je liever je eigen namen voor dataframe rijen en -kolommen opgeeft, kunnen deze worden toegewezen met behulp van de functies **rownames()** en **colnames()** - net als bij matrices. Je kunt ook een argument **row.range=** opnemen in de aanroep van de functie **data.frame()** om een lijst met namen voor de dataframe rijen op te geven.

Individuele waarden worden geadresseerd in een dataframe met behulp van de juiste indexnummers van de rij en kolom. Nieuwe waarden kunnen ook aan individuele cellen worden toegewezen met behulp van hun indexnummer van rij en kolom, maar let op het gegevenstype van die cel.

Om te bevestigen dat een structuur inderdaad een dataframe is, biedt R een ingebouwde functie **is.data.frame()** die alleen **TRUE** retourneert als het opgegeven argument een dataframe-object is.



Het is handig om een waarde te zoeken met in een gegevensframe met behulp van de ingebouwde functie **which()** om een logisch argument op te geven en optioneel kunt u een argument **arr.ind=TRUE** opnemen om het rij- en kolom-indexnummer van elke cel te retourneren met daarin de gezochte waarde.

- 1 Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.
- 2 Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.

```
1 #Maak drie vectoren met verschillende gegevenstypewaarden.
2 bools <- c(TRUE, FALSE, TRUE)
3 chars <- LETTERS[1:3]
4 nums <- 1:3
5
6 #Maak een dataframe dat de vectorwaarden opslaat, voer het
7 #dataframe uit en bevestig de structuur ervan.
8 frame <- data.frame(bools, chars, nums)
9 print(frame)
10 cat("\nData Frame?: ", is.data.frame(frame), "\n\n")
11
12 #Geef de rijen en kolommen van het dataframe een naam
13 rownames(frame) <- c("Tier 1:", "Tier 2:", "Tier 3: ")
14 colnames(frame) <- c("Logical", "Alphabetical", "Mumerical")
15
16 #wijs een nieuwe waarde toe aan één cel en voer het gegevensframe
17 #nogmaals uit om de benoemde rijen en kolommen te zien.
18 frame[2,2] <- "A"
19 print(frame)
20
21 #Zoek in alle cellen naar een specifieke waarde en identificeer
22 #de locatie van cellen die die waarde wel bevatten.
23 cat("Search for 'A'...\n")
24 print(which(frame == "A", arr.ind=TRUE))
25
26
```



Als alternatief kunnen de rij-namen worden toegewezen aan een argument **row.names=** in de aanroep van de functie **data.frames()**.

- 3 Sla het **R Script** op met de naam **FirstDataframe.R** en klik vervolgens op de knop **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```
Console Terminal x Jobs x
~/
> source('E:/IT Concepten/IT Concepten - New Site/1
essen/rstudio/sourcecode/FirstDataFrame2.R')
  bools chars nums
1  TRUE    A    1
2 FALSE    B    2
3  TRUE    C    3

Data Frame?: TRUE

      Logical Alphabetical Numerical
Tier 1:   TRUE           A           1
Tier 2:  FALSE           A           2
Tier 3:   TRUE           C           3
Search for 'A'...
      row col
Tier 1:   1  2
Tier 2:   2  2
> |
```

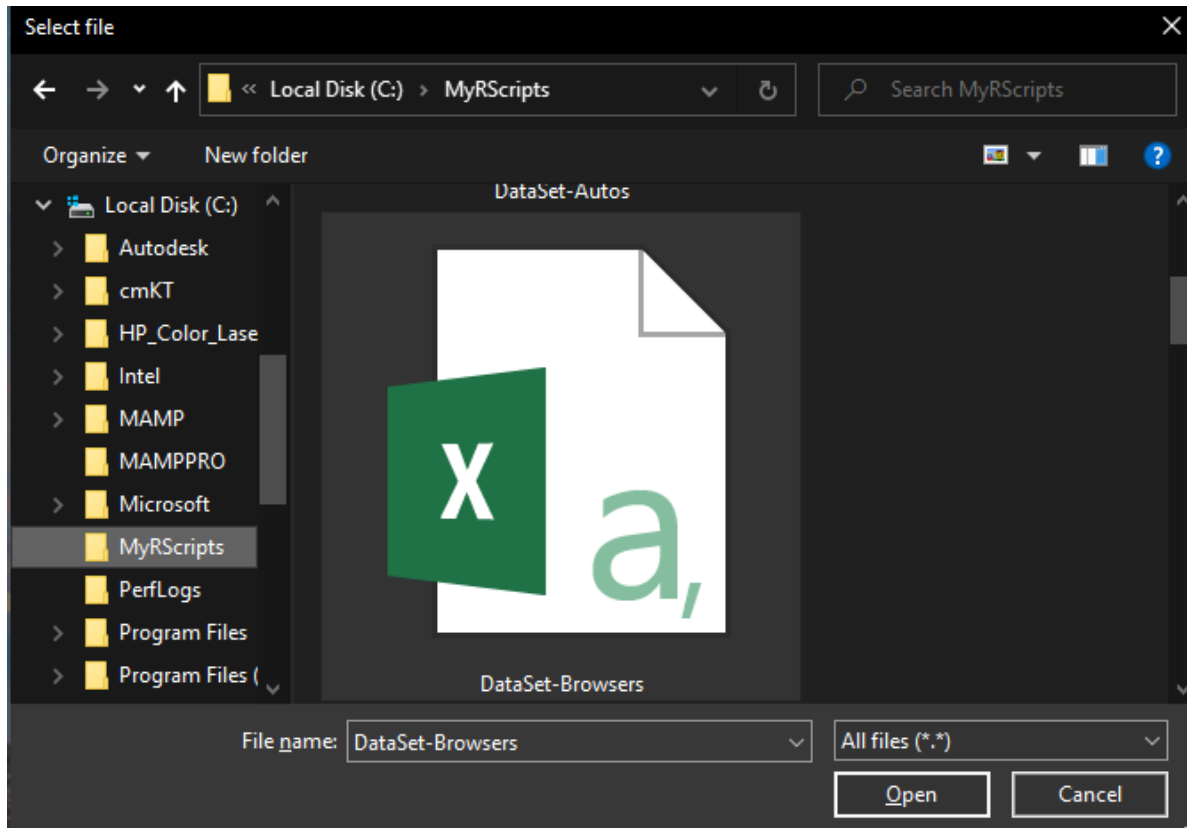
Importing data sets

Verzamelingen van gegevens in tabelvorm worden vaak opgeslagen als een "**data set**" (gegevensset) in een bestand met door komma's gescheiden waarden (CSV = Comma-Separated Values). Deze datasets kunnen eenvoudig worden geïmporteerd in **RStudio** en hun gegevens kunnen worden gekopieerd naar een dataframe met behulp van de ingebouwde functie **read.csv()**. Deze functie vereist eenvoudigweg het CSV-bestands-pad als argument.

Het pad naar een CSV-bestand kan worden opgegeven voor de functie **read.csv()** door de functie **file.choose()** als argument op te geven. Hierdoor wordt een dialoogvenster "Bestand selecteren" geopend waarmee u naar de locatie van het CSV-bestand kunt bladeren. Eenmaal geselecteerd, wordt het pad geleverd aan de functie **read.csv()** zodat de gegevens naar een gegevensframe kunnen worden gekopieerd. Het instructie **frame <- read.csv(file.choose())** biedt bijvoorbeeld een dialoogvenster om een CSV-bestand te selecteren waarvan de gegevens vervolgens worden gekopieerd naar een gegevensframe met de naam "*frame*".

Als alternatief kan het volledige pad van een CSV-bestand worden opgegeven voor de functie **read.csv()** als argument. Deze instructie

geeft bijvoorbeeld het volledige pad naar het geselecteerde CSV-bestand dat hieronder wordt geïllustreerd:



```
frame <- read.csv("C:/MyRScripts/DataSet-Browsers.csv")
```

Volledige pad-namen kunnen lang zijn, dus het is handig om de locatie van het CSV-bestand in te stellen als de werkdirectory van **RStudio** door de directory op te geven als argument voor de functie **setwd()**. De functie **read.csv()** vereist dan alleen de CSV-bestandsnaam als argument. U kunt de huidige werkdirectory op elk moment vinden door de ingebouwde **getwd()**-functie aan te roepen.



JE kunt de standaard werkmap in **RStudio** instellen door op **Tools, Global Options** en de optie **General** te klikken.



Paden kunnen de voorwaartse slash / teken of escaped backslash gebruiken, zoals **C:\\MyRScripts\\DataSets-Browser.csv**

1

Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.

- 2 Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.

```

1 #Stel de werkmapp in en bevestig vervolgens het padadres.
2 setwd("c:/MyRScripts")
3 cat("working Directory:", getwd(), "\n\n")
4
5 #Maak een dataframe dat gegevens importeert uit een
6 #CSV-bestand dat zich in de werkmapp bevindt.
7 frame <-read.csv("DataSet-Browsers.CSV")
8
9 #voer de gegevensframerijen en -kolommen uit.
10 print(frame)
11

```

- 3 Sla het **R Script** op met de naam **ImportData.R** en klik vervolgens op de knop **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```

Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/lessen/rstudio/sourcecode/ImportData2.R')
Working Directory: c:/MyRScripts

      Web.Browser.Version Percentage.Market.Share
1      Chrome 61.0          35.16
2 Microsoft Internet Explorer 11.0          12.52
3      Chrome 55.0           7.40
4      Firefox 56           6.53
5      Microsoft Edge 14           3.60
6      Chrome 49.0           2.26
7      Chrome 60.0           2.19
8      Chrome 45.0           1.86
9      Chrome 62.0           1.76
10     Safari 11.0           1.70
11     Chrome 50.0           1.59
12     Microsoft Edge 15           1.43
13     Chrome 58.0           1.25
14     Firefox 55           1.24
15 Microsoft Internet Explorer 8.0           1.12
16     Firefox 52           1.06
17     Chrome 53.0           0.93
18     Safari 10.0           0.85
19     Opera 48              0.85
20     Firefox 53           0.83
21 Microsoft Internet Explorer 9.0           0.81
22     Firefox 40           0.76
23     Safari 10.1           0.66
24     Microsoft Edge 13           0.65
25     Chrome 47.0           0.64
26     Chrome 59.0           0.56
27 Microsoft Internet Explorer 10.0           0.54
28     Safari 8.0            0.52
29     Chrome 56.0           0.49
30     Firefox 57           0.47
31     Chrome 57.0           0.47
32     Chrome 43.0           0.41
33     Microsoft Edge 12           0.33
34     Chrome 60.4           0.28
35     Safari 9.1            0.26
36     Firefox 50           0.21
37     Firefox 54           0.20
38     Firefox 8.0           0.19
39     Firefox 47           0.19
40     Chrome 54.0           0.17
>

```



Hier worden de kolomnamen geleverd door de eerste regel van het CSV-bestand, maar de rijen worden automatisch genummerd door de R-interpretter. Merk op dat er drie gevallen van identiek marktaandeel zijn.

Examining data frames

De programmeertaal **R** biedt een aantal functies voor het onderzoeken van data frame structuren:

- **nrow()**
 - retourneert een geheel getal dat het totale aantal rijen is binnen het gegevensframe dat is opgegeven als argument.
- **ncol()**
 - retourneert een geheel getal dat het totale aantal kolommen is binnen het gegevensframe dat is opgegeven als argument.
- **head()**
 - retourneert standaard de bovenste zes rijen van het dataframe dat is opgegeven als zijn argument, plus de kolom- en rijnamen. Optioneel kan een argument **n=** worden toegevoegd om aan te geven hoeveel rijen moeten worden geretourneerd. **n=3** retourneert bijvoorbeeld de bovenste drie rijen.
- **tail()**
 - retourneert standaard de onderste zes rijen van het dataframe dat is opgegeven als argument, plus kolom- en rijnamen. Een argument **n=** kan worden toegevoegd om het aantal te retourneren rijen op te geven. **n=3** retourneert bijvoorbeeld de onderste drie rijen.
- **str()**
 - geeft een overzicht weer van de structuur van het dataframe dat als argument is opgegeven - met een lijst van het totale aantal objecten en variabelen dat het bevat, samen met het aantal unieke factoren (**factors**, *zijn unieke waarden*) in elke kolom. Intern worden de factoren numeriek gerangschikt in aflopende volgorde, waarbij de hoogste rangorde niveau 1 is.
- **summary()**
 - geeft een samenvatting weer van de inhoud van het dataframe dat is opgegeven als argument. Standaard worden zes niveaus van factoren weergegeven, maar een optioneel **maxsum=**-argument kan worden toegevoegd om te specificeren hoeveel niveaus moeten worden weergegeven. Voor kolommen die numerieke gegevens bevatten, biedt de samenvatting statistieken:
 - **Minimum** - het laagste getal in de kolom.
 - **1e kwartiel** (1st Quartile) - de middelpuntwaarde tussen het minimumaantal en de mediaanwaarde.
 - **Mediaan** (Median) - de middelste waarde van de kolom.
 - **Gemiddelde** (Mean) - de gemiddelde waarde van de getallen in de kolom.

- **3e kwartiel** (3rd Quartile) - de middelpuntwaarde tussen de mediaanwaarde en het maximale aantal.
- **Maximum** - het hoogste getal in de kolom.

1

Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.


2

Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.

```
ExamineData.R* x
Source on Save
Run
Source

1 #Stel de werkmapijn
2 setwd("c:/MyRscripts")
3
4 #Maak een gegevensframe dat gegevens importeert uit een
5 #CSV-bestand dat zich in de werkmapijn bevindt.
6 frame <- read.csv("DataSet-Browsers.csv")
7
8 #Geef het totale aantal rijen en kolommen weer.
9 cat("Rows: ", nrow(frame), "\tColumns: ", ncol(frame))
10
11 #Voer de eerste en laatste drie rijen en kolommen uit
12 cat("\nHead...\n")
13 print(head(frame, n=3))
14 cat("\nTail...\n")
15 print(tail(frame, n=3))
16
17 #Toon de structuur en samenvatting
18 cat("\nStructure...\n")
19 print(str(frame))
20 cat("\nSummary...\n")
21 print(summary(frame))
22
```

3

Sla het **R Script** op met de naam **ExaminData.R** en klik vervolgens op de knop  **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.



Het is handig om te controleren of het totale aantal geïmporteerde rijen en kolommen overeenkomt met wat wordt verwacht van de bron-gegevens-set.



In R-programmering verwijst de functie `str()` naar "structuur", in tegenstelling tot andere programmeertalen waarin het verwijst naar "string".

```

Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/lessen/rstudio/
sourcecode/ExamineData2.R')
Rows: 40      Columns: 2
Head...
      web.Browser.Version PerCentage.Market.Share
1          Chrome 61.0          35.16
2 Microsoft Internet Explorer 11.0          12.52
3          Chrome 55.0          7.40

Tail...
      web.Browser.Version PerCentage.Market.Share
38          Firefox 8.0          0.19
39          Firefox 47          0.19
40          Chrome 54.0          0.17

Structure...
'data.frame':  40 obs. of  2 variables:
 $ web.Browser.Version  : chr  "Chrome 61.0" "Microsoft Interne
t Explorer 11.0" "chrome 55.0" "Firefox 56" ...
 $ PerCentage.Market.Share: num  35.16 12.52 7.4 6.53 3.6 ...
NULL

Summary...
web.Browser.Version PerCentage.Market.Share
Length:40          Min.   : 0.170
Class :character  1st Qu.: 0.470
Mode  :character  Median : 0.820
                  Mean  : 2.373
                  3rd Qu.: 1.617
                  Max.   :35.160
> |

```



Hier bevat de eerste kolom gegevens van het gegevenstype karakter, terwijl de tweede kolom gegevens van het gegevenstype numeriek dubbel bevat - dus statistieken worden voor die kolom in de samenvatting gegeven.

Addressing frame data

Er zijn een aantal manieren om gegevens in een dataframe aan te pakken. Net als bij matrices kunnen cellen worden geadresseerd door hun rij- en kolom index nummer tussen [] vierkante haken te vermelden. Als alternatief kunnen cellen worden geadresseerd door hun rij- en kolomnaam tussen [] vierkante haken te vermelden. Bovendien kunnen dataframes **de \$ dollar-operator** gebruiken, zodat een kolom kan worden geadresseerd met deze handige syntaxis:

data.frame.name\$column.name


Evenzo kunnen cellen worden geadresseerd door [] vierkante haken in de bovenstaande syntaxis waarin rij-indexnummers moeten worden opgegeven. De factorniveaus van een kolom kunnen worden

geadresseerd door de kolom op te geven als argument voor de ingebouwde functie **levels()** met dezelfde syntaxis:

- 1 Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.
- 2 Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.

```
AddressData.R x
Source on Save
Run
Source

1 #Stel de werkmapp in
2 setwd("C:/MyRscripts")
3
4 #Maak een gegevensframe dat gegevens importeert uit een
5 #CSV-bestand dat zich in de werkmapp bevindt.
6 frame <- read.csv("DataSet-Browsers.csv")
7
8 #Voer de eerste drie rijen en kolommen uit.
9 cat("\nHead...\n")
10 print(head(frame, n=3))
11
12 #Adresseer een enkele cel per rij- en kolomindexnummer.
13 data <- frame[ 1, 2]
14 cat("\nRow #1, Column #2:", data, "\n")
15
16 #Adresseer een enkele cel op rij-indexnummer en kolomnaam.
17 data <- frame[ 2, "Percentage.Market.Share"]
18 cat("\nRow #2, Column #2:", data, "\n")
19
20 #Adresseer een enkele cel op kolomnaam en rijindex.
21 data <- frame$Percentage.Market.Share[3]
22 cat("\nRow #3, Column #2:", data, "\n")
23
24 #Adresseer een hele kolom op naam om alle niveaus uit
25 #te voeren.
26 print(levels(factor(frame$web.Browser.Version)))
27
```

- 3 Sla het **R Script** op met de naam **AddressData.R** en klik vervolgens op de knop  **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```
Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/lessen/rstudio/sourcocode/AddressData3.R')

Head...
      web.Browser.Version Percentage.Market.Share
1          Chrome 61.0          35.16
2 Microsoft Internet Explorer 11.0      12.52
3          Chrome 55.0          7.40

Row #1, Column #2: 35.16
Row #2, Column #2: 12.52
Row #3, Column #2: 7.4
[1] "Chrome 43.0"
[2] "Chrome 45.0"
[3] "Chrome 47.0"
[4] "Chrome 49.0"
[5] "Chrome 50.0"
[6] "Chrome 53.0"
[7] "Chrome 54.0"
[8] "Chrome 55.0"
[9] "Chrome 56.0"
[10] "Chrome 57.0"
[11] "Chrome 58.0"
[12] "Chrome 59.0"
[13] "Chrome 60.0"
[14] "Chrome 60.4"
[15] "Chrome 61.0"
[16] "Chrome 62.0"
[17] "Firefox 40"
[18] "Firefox 47"
[19] "Firefox 50"
[20] "Firefox 52"
[21] "Firefox 53"
[22] "Firefox 54"
[23] "Firefox 55"
[24] "Firefox 56"
[25] "Firefox 57"
[26] "Firefox 8.0"
[27] "Microsoft Edge 12"
[28] "Microsoft Edge 13"
[29] "Microsoft Edge 14"
[30] "Microsoft Edge 15"
[31] "Microsoft Internet Explorer 10.0"
[32] "Microsoft Internet Explorer 11.0"
[33] "Microsoft Internet Explorer 8.0"
[34] "Microsoft Internet Explorer 9.0"
[35] "Opera 48"
[36] "Safari 10.0"
[37] "Safari 10.1"
[38] "Safari 11.0"
[39] "Safari 8.0"
[40] "Safari 9.1"
> |
```



Merk op hoe de niveaus automatisch worden aangemaakt in aflopende alfabetische en numerieke volgorde.

Extracting frame subsets

Subsets kunnen worden geëxtraheerd uit dataframes door hun rij- en kolomindex nummers tussen [] vierkante haken te vermelden - op dezelfde manier als subsets worden geëxtraheerd uit matrices. Gebruik bijvoorbeeld **frame[1:3, 2]** om een subset te extraheren die alleen de tweede kolom van de eerste drie rijen gegevensframe bevat, of **frame[1:3,]** om een subset te extraheren die alle kolommen van de eerste drie rijen bevat.

Als alternatief kan een subset van specifieke cellen worden geëxtraheerd met behulp van de bekende **c()**-functie. Als u bijvoorbeeld **frame[c(1, 3, 5), 2]** gebruikt om een subset te extraheren die alleen de tweede kolom van de eerste, derde en vijfde rij van een gegevensframe bevat, of **frame[c(1, 3, 5),]** om een subset te extraheren die alle kolommen van die specifieke rijen bevat.

Wanneer een subset een of meer rijen of meer dan één kolom extraheert uit een dataframe, worden de gegevens geretourneerd in een dataframe-object. Omgekeerd, wanneer een subset slechts één kolom extraheert, worden de gegevens standaard geretourneerd in een vectorobject. Het standaardgedrag kan worden omzeild door een laatste argument **drop=FALSE** op te nemen tussen de [] vierkante haken die de kolom specificeren die moet worden geëxtraheerd - zodat de gegevens nu worden geretourneerd in een vectorstructuur:

1 Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.

2 Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.



Vergeet de extra spatie en komma binnen de vierkante haken niet om een hele rij of een hele kolom aan te duiden. Regels 16, 30 en 41.




```
SubsetData.R x
Source on Save
Run
Source

1 #Stel de werkmapi in
2 setwd("C:/MyRScripts")
3
4 #Maak een gegevensframe dat gegevens importeert uit een
5 #CSV-bestand dat zich in de werkmapi bevindt.
6 frame <- read.csv("DataSet-Browsers.csv")
7
8 #Extraheer een subset met gegevens in alle kolommen van
9 #vier specifieke rijen en voer vervolgens de subset-
10 #gegevens uit.
11 edge <- frame[c(33,24,5, 12), ]
12 print(edge)
13
14 #Extraheer uit de oorspronkelijke subset een tweede
15 #subset met gegevens in alle kolommen van één enkele rij.
16 edge.row <- edge[1, ]
17
18 #voer de gegevens van de tweede subset uit om te zien
19 #dat de kolomkoppen en rijnummering behouden blijven.
20 cat("\nRow...\n")
21 print(edge.row)
22
23 #Bevestig de gegevensstructuur van de tweede subset.
24 cat("Data Frame?:", is.data.frame(edge.row))
25
26 #Extraheer uit de oorspronkelijke subset een derde
27 #subset die alle gegevens in één enkele kolom bevat en
28 #voer vervolgens de gegevens uit om te zien dat
29 #kolomkoppen en rijnummers niet behouden blijven.
30 edge.col <- edge[ ,2]
31 cat("\n\nColumn...\n")
32 print(edge.col)
33
34 #Bevestig de gegevensstructuur van de derde subset
35 cat("Data Frame?: ", is.data.frame(edge.col))
36 cat("\tVector?: ", is.vector(edge.col))
37
38 #Maak de derde subset opnieuw om het standaardgedrag
39 #te negeren en voer vervolgens de gegevens uit om te
40 #zien dat kolomkoppen en rijnummers nu behouden blijven.
41 edge.col <- edge[ , 2, drop=FALSE]
42 cat("\n\nColumn...\n")
43 print(edge.col)
44
45 #Bevestig de gegevensstructuur van de opnieuw gemaakte
46 #derde subset.
47 cat("Data Frame?:", is.data.frame(edge.col))
48 cat("\tVector?: ", is.vector(edge.col))
49
```



3

Sla het **R Script** op met de naam **SubsetData.R** en klik vervolgens op de knop  **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```
Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/lessen/rstudio/
sourcecode/SubsetData2.R')
  Web.Browser.Version PerCentage.Market.Share
33 Microsoft Edge 12 0.33
24 Microsoft Edge 13 0.65
5 Microsoft Edge 14 3.60
12 Microsoft Edge 15 1.43

Row...
  Web.Browser.Version PerCentage.Market.Share
33 Microsoft Edge 12 0.33
Data Frame?: TRUE

Column...
[1] 0.33 0.65 3.60 1.43
Data Frame?: FALSE Vector?: TRUE

Column...
  PerCentage.Market.Share
33 0.33
24 0.65
5 3.60
12 1.43
Data Frame?: TRUE Vector?: FALSE
> |
```

Changing frame columns

De kolomnamen van gegevensframes die zijn gemaakt op basis van gegevens die zijn geïmporteerd uit een CSV-bestand, worden meestal overgenomen van de koptekstnamen die zijn opgegeven op de eerste regel van het CSV-bestand. Als u de headernamen liever negeert, kunt u een **header=FALSE**-argument toevoegen in de aanroep van de **read.csv()**-functie.

Kolomnamen kunnen worden toegevoegd aan naamloze kolommen, of bestaande kolomnamen kunnen worden gewijzigd door kolomnamen toe te wijzen aan de functie **colnames()**, zoals bij matrices.

Een nieuwe kolom kan aan een gegevensframe worden toegevoegd door cel-waarden als nieuwe kolomnaam toe te wijzen, zoals deze:

```
frame["New Column"] <- 1:10
```

Als alternatief kan de \$ dollar-operator worden gebruikt om een nieuwe kolom toe te voegen, maar een kolomnaam die spaties bevat, moet worden ingesloten tussen ` **backtick**-tekens, zoals deze:

```
frame$`New Column` <- 1:10
```



De ` Backtick-toets bevindt zich meestal linksboven op het toetsenbord, naast de cijfertoets 1.

Een bestaande kolom kan eenvoudig uit een gegevensframe worden verwijderd door er een NULL-waarde aan toe te kennen:

```
frame$Existing.Column <- NULL
```

Waarden die aan een nieuwe kolom zijn toegewezen, kunnen worden gekopieerd uit een bestaande kolom en worden gemanipuleerd om nieuwe gegevenswaarden te verschaffen. Om bijvoorbeeld twee kolommen in een derde kolom te tellen:

```
frame$Total <- frame$Price + frame$Tax
```

Numerieke gegevens kunnen worden geconverteerd naar het tekengegevenstype voor tekenreeksmanipulatie door de numerieke gegevenswaarde op te geven als argument voor de ingebouwde functie **as.character()**. Het kan vervolgens worden samengevoegd met andere tekengegevens met behulp van de functie **paste()**. Standaard voegt de functie `paste()` automatisch spatietekens toe als scheidingstekens tussen aaneengeschakelde tekenreeksen. Dit gedrag kan echter worden overschreven door een alternatief scheidingsteken op te geven voor een optioneel **sep=**-argument voor de functie **paste()**. Als je liever geen scheidingsteken hebt tussen aaneengeschakelde tekenreeksen, kan dit argument helemaal geen tekens bevatten.



Rekenkundige en wiskundige bewerkingen kunnen alleen worden uitgevoerd op numerieke gegevenswaarden, zoals dubbele of integer gegevenstypewaarden.


1

Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.

2

Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.

```
ColumnData.R* x
Source on Save Run Source
1 #Stel de werkmapi in
2 setwd("C:/MyRScripts")
3
4 #Maak een gegevensframe dat gegevens importeert uit een
5 #CSV-bestand dat zich in de werkmapi bevindt.
6 frame <- read.csv("DataSet-Browsers.csv")
7
8 #Maak een functie om een titel en de eerste twee regels
9 #van een dataframe uit te voeren - inclusief rij- en
10 #kolomnamen.
11 display <- function(title){
12   cat("\n", title, "...\n")
13   print(head(frame, n=2))
14 }
15
16 #voer twee gegevensframerijen en -kolommen uit, hernoem
17 #vervolgens de bestaande kolommen en voer de wijzigingen
18 #uit.
19 display("Original columns")
20 colnames(frame) <- c("web.Browser", "PerCentage")
21 display("Renamed columns")
22
23 #Kopieer bestaande numerieke kolomgegevens naar een nieuwe
24 #kolom - geconverteerd naar tekengegevens en aaneengeschaeld.
25 frame$Market.Share <-
26   paste(as.character(frame$PerCentage), "%", sep="")
27
28 #Verwijder de kolom met numerieke gegevens en voer de
29 display(title) igen nogmaals uit.
30 frame$PerCentage <- NULL
31 display("Switch columns")
32
```

3 Sla het **R Script** op met de naam **ColumnData.R** en klik vervolgens op de knop  **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.



Het verwijderen van kolommen heeft invloed op hun indexnummer. Als u bijvoorbeeld de eerste kolom verwijdert, betekent dit dat de tweede kolom index nummer één aanneemt.

```

Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/lessen/rstudio/sourcecode/ColumnData2.R')

Original Columns ...
      Web.Browser.Version PerCentage.Market.Share
1          Chrome 61.0          35.16
2 Microsoft Internet Explorer 11.0          12.52

Renamed Columns ...
      Web.Browser PerCentage
1          Chrome 61.0          35.16
2 Microsoft Internet Explorer 11.0          12.52

Switch Columns ...
      Web.Browser Market.Share
1          Chrome 61.0          35.16%
2 Microsoft Internet Explorer 11.0          12.52%
> |

```

Filtering data frames


De gegevens in de cellen van een gegevensframe kunnen worden gefilterd door een voorwaardelijke test uit te voeren op de waarde in elke cel, om slechts één gegevens te selecteren uit cellen waarvan de voorwaardelijke test **TRUE** retourneert.

Numerieke gegevens kunnen worden gefilterd door de waarde in elke cel te vergelijken met een waarde die is opgegeven in een voorwaardelijke test. Er kan bijvoorbeeld een filter worden gemaakt voor een kolom die numerieke gegevens bevat, om de pariteit (gelijkheid/overeenkomst) van gegevens in elke cel te onderzoeken en alleen even waarden te selecteren met dit instructie **filter <- frame[,1] %% 2 == 0**. Het filter is gemaakt als een vector van Booleaanse waarden waarin even waarden worden weergegeven door **TRUE** en oneven waarden door **FALSE**. Als u het filter opneemt tussen de **[]** vierkante haken die een gebied van het gegevensframe specificeren, worden alleen cellen geselecteerd waarvan de gegevens de voorwaardelijke test doorstaan:

```

Console Terminal x Jobs x
~/
> frame <- data.frame(A=1:5,B=NA)
> filter <- frame[, 1] %% 2 == 0
> print(t(frame))
  [,1] [,2] [,3] [,4] [,5]
A    1    2    3    4    5
B   NA   NA   NA   NA   NA
> print(filter)
[1] FALSE TRUE FALSE TRUE FALSE
> print(frame$A[filter])
[1] 2 4
> |

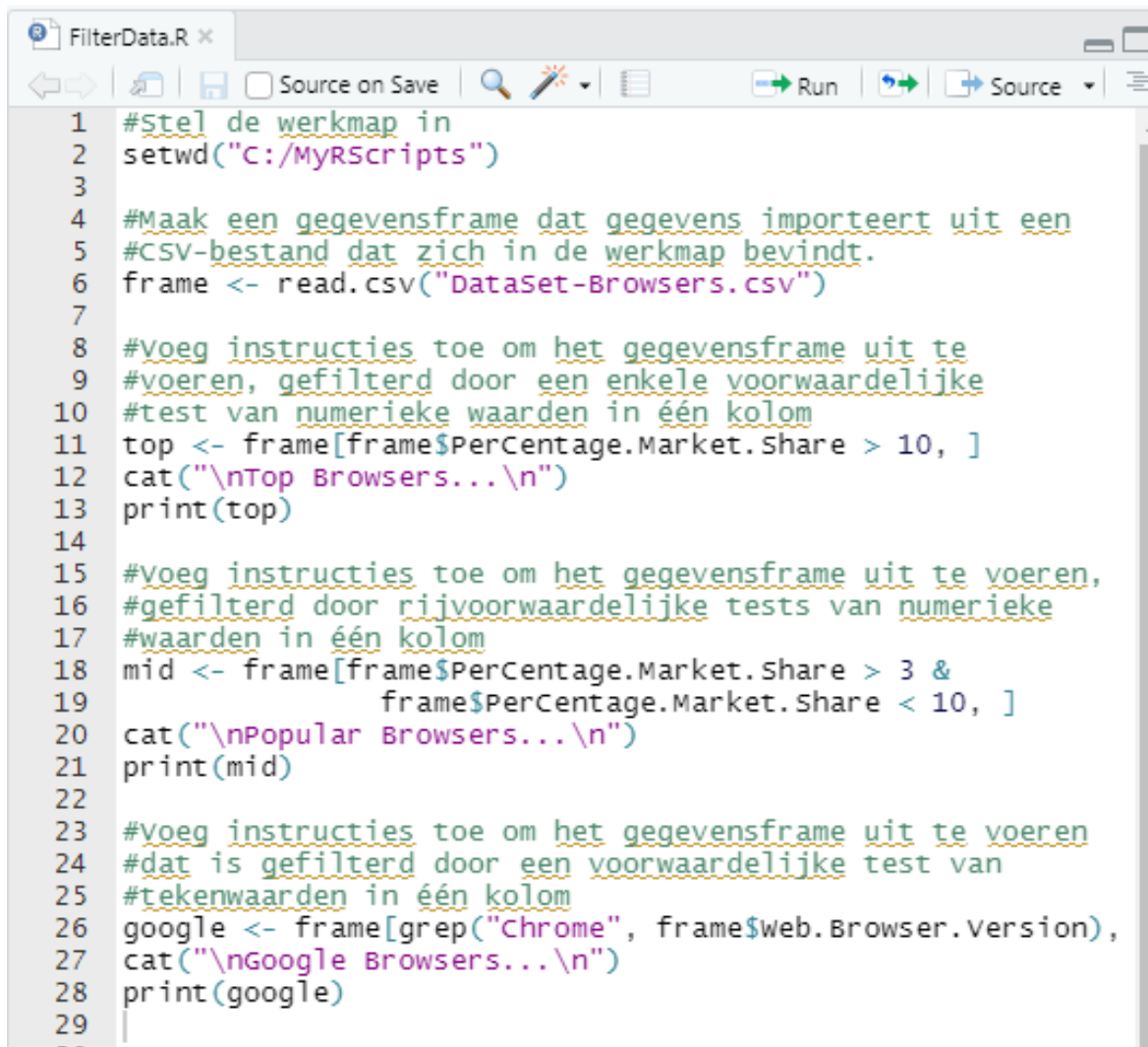
```

 De functie **t()** wordt gebruikt om het gegevensframe te transponeren, zodat de kolommen ervan als rijen worden weergegeven.

Tekengegevens kunnen worden gefilterd door de waarde in elke cel te vergelijken met een waarde die is opgegeven als het eerste argument voor de ingebouwde **grep()**-functie. Deze functie vereist ook een tweede argument om het gebied van het gegevensframe op te geven waarvan de cel waarden moeten worden vergeleken.

Het voorwaardelijke testresultaat van een filter kan worden toegewezen aan een variabele, zoals in de bovenstaande schermafbeelding, of de voorwaardelijke test kan worden gemaakt tussen de [] vierkante haken die een gebied van het gegevensframe specificeren:

- 1 Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.
- 2 Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.



```
1 #Stel de werkmapp in
2 setwd("C:/MyRScripts")
3
4 #Maak een gegevensframe dat gegevens importeert uit een
5 #CSV-bestand dat zich in de werkmapp bevindt.
6 frame <- read.csv("DataSet-Browsers.csv")
7
8 #Voeg instructies toe om het gegevensframe uit te
9 #voeren, gefilterd door een enkele voorwaardelijke
10 #test van numerieke waarden in één kolom
11 top <- frame[frame$PerCentage.Market.Share > 10, ]
12 cat("\nTop Browsers...\n")
13 print(top)
14
15 #Voeg instructies toe om het gegevensframe uit te voeren,
16 #gefilterd door rijvoorwaardelijke tests van numerieke
17 #waarden in één kolom
18 mid <- frame[frame$PerCentage.Market.Share > 3 &
19               frame$PerCentage.Market.Share < 10, ]
20 cat("\nPopular Browsers...\n")
21 print(mid)
22
23 #voeg instructies toe om het gegevensframe uit te voeren
24 #dat is gefilterd door een voorwaardelijke test van
25 #tekenwaarden in één kolom
26 google <- frame[grep("Chrome", frame$web.Browser.Version),
27                 cat("\nGoogle Browsers...\n")
28                 print(google)
29
```



De functie **grep()** kan optioneel een argument **ignore.case=TRUE** bevatten als hoofdlettergevoeligheid niet vereist is.

3

Sla het **R Script** op met de naam **ColumnData.R** en klik vervolgens op de knop **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```
Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/lessen/rstudio/s
ourcecode/FilterData2.R')

Top Browsers...
      Web.Browser.Version PerCentage.Market.Share
1          Chrome 61.0          35.16
2 Microsoft Internet Explorer 11.0          12.52

Popular Browsers...
      Web.Browser.Version PerCentage.Market.Share
3          Chrome 55.0          7.40
4          Firefox 56          6.53
5 Microsoft Edge 14          3.60

Google Browsers...
      Web.Browser.Version PerCentage.Market.Share
1          Chrome 61.0          35.16
3          Chrome 55.0          7.40
6          Chrome 49.0          2.26
7          chrome 60.0          2.19
8          Chrome 45.0          1.86
9          chrome 62.0          1.76
11         Chrome 50.0          1.59
13         Chrome 58.0          1.25
17         Chrome 53.0          0.93
25         Chrome 47.0          0.64
26         Chrome 59.0          0.56
29         Chrome 56.0          0.49
31         Chrome 57.0          0.47
32         Chrome 43.0          0.41
34         chrome 60.4          0.28
40         Chrome 54.0          0.17
> |
```

De functie **grep()** wordt hier gebruikt om alle rijen te selecteren waarvan de eerste kolom cel de tekenreeks "Chrome" bevat.

Merging data frames

Gegevens uit twee gegevensframes kunnen worden gecombineerd in een enkel gegevensframe door een gemeenschappelijke kolomnaam of een gemeenschappelijke rij-naam met behulp van de ingebouwde functie **merge()**. Deze functie vereist de namen van de twee bron data frames als de eerste twee argumenten. Het vereist ook dat de namen van het kolom- of rijveld met gemeenschappelijke waarden worden opgegeven voor de argumenten **by.x=** en **by.y=**.



Als twee dataframes worden samengevoegd door gemeenschappelijke velden die identieke namen hebben, kunnen de argumenten **by.x=** en **by.y=** worden vervangen door één enkel argument **by=** dat de naam van het gemeenschappelijke veld vermeldt. Bijvoorbeeld **by="State"**.

Als je er de voorkeur aan geeft alle velden te behouden, ongeacht of ze precies overeenkomen, kunt je een **all=TRUE** argument opnemen in de aanroep van de **merge()**-functie. In dit geval worden cellen die niet direct overeenkomen in elk bron data frame, gevuld met de speciale **NA** (*Not Available = niet beschikbaar*) constante in het samengevoegde data frame.

De twee algemene velden van de bron data frames verschijnen als een enkel veld in de samengevoegde data frame, maar andere velden die gegevens bevatten die

gemeenschappelijk zijn voor beide bron data frames, verschijnen in afzonderlijke velden. Dit betekent dat het samengevoegde data frame rijen of kolommen met geduplicateerde gegevens kan bevatten. Door een **NULL**-waarde toe te wijzen aan onnodige velden, worden ze verwijderd, zodat alleen unieke gegevens in het samengevoegde data frame overblijven.

Extra kolommen kunnen zoals gebruikelijk aan een samengevoegd data frame worden toegevoegd en kunnen gegevens bevatten die zijn berekend op basis van de gegevens die uit elk bron data frame zijn geleverd:

1

Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.

2


Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.



In de onderstaande broncode hebben de kolomvelden met staats-afkortingen verschillende namen en verschijnen op verschillende posities binnen de bron data frames.

```
MergeData.R x
Source on Save
Run
Source

1 #Stel de werkmapp in
2 setwd("C:/MyRscripts")
3
4 #Maak twee data frames van gegevensimport uit csv-bestanden
5 #in de werkmapp
6 high.temp <- read.csv("DataSet-HighTemps.csv")
7 low.temp <- read.csv("DataSet-LowTemps.csv")
8
9 #Maak een functie om een titel en een dataframe uit te voeren
10 #wanneer aangeropen
11 display <- function(frame){
12   cat("\nAnnual Temperatures (°C)... \n")
13   print(frame)
14 }
15
16 #Voeg een instructie toe om elk brongegevensframe uit te voeren
17 display(high.temp)
18 display(low.temp)
19
20 #Voeg een instructie toe om de twee gegevensframes samen te
21 #voegen op basis van gemeenschappelijke kolomwaarden
22 avg.temp <- merge(high.temp, low.temp,
23                  by.x="state", by.y="State.Code")
24 display(avg.temp)
25
26 #Voeg een instructie toe om een kolom te verwijderen die dubbele
27 #gegevens bevat
28 avg.temp$capital <- NULL
29
30 #Voeg een nieuwe kolom toe met gegevens die zijn berekend op
31 #basis van de gegevens die zijn geleverd door de twee bron
32 #data frames en voer de gecombineerde gegevens nogmaals uit
33 avg.temp$Average <- (avg.temp$High + avg.temp$Low)/2
34 display(avg.temp)
35
```

- 3 Sla het **R Script** op met de naam **MergeData.R** en klik vervolgens op de knop  **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```

Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/les
sen/rstudio/sourcecode/MergeData2.R')

Annual Temperatures (°C)...
  City State High
1 Albany NY 14.3
2 Sacramento CA 23.1
3 Austin TX 26.5

Annual Temperatures (°C)...
  State.Code Capital Low
1 NY Albany 3.7
2 CA Sacramento 9.1
3 TX Austin 15.0

Annual Temperatures (°C)...
  State City High Capital Low
1 CA Sacramento 23.1 Sacramento 9.1
2 NY Albany 14.3 Albany 3.7
3 TX Austin 26.5 Austin 15.0

Annual Temperatures (°C)...
  State City High Low Average
1 CA Sacramento 23.1 9.1 16.10
2 NY Albany 14.3 3.7 9.00
3 TX Austin 26.5 15.0 20.75
> |

```



De stadnamen worden gedupliceerd in het samengevoegde data frame, zodat één kolom niet nodig is en kan worden verwijderd.

Adjusting factors

De structuur van gegevens die vanuit een data frame naar een vectorvariabele zijn gekopieerd, blijft behouden, zodat de **factors** en **levels** behouden blijven. Vectoren bieden niet automatisch factorcategorieën en gerangschikte levels, maar de naam van een vectorvariabele kan worden opgegeven als argument voor de ingebouwde **factor()**-functie om **factors** en **levels** te creëren.

De functie **factor()** rangschikt standaard de factor levels (niveaus) in alfabetische en numerieke volgorde. Als de gegevenswaarden van bijvoorbeeld "B", "C2", "A" en "C1" worden gegeven, worden de factorniveaus gerangschikt als "A" (1), "B" (2), "C1" (3), en "C2" (4).

U kunt optioneel een rangordevoorkeur opgeven door een argument **levels=** op te nemen in de aanroep van de functie **factor()**. Aan dit argument moet een vector van de data values (gegevenswaarden) worden toegewezen in de aflopende volgorde van het rangordeniveau van uw voorkeur, **levels=c("C1", "C2", "A", "B")** zorgt er bijvoorbeeld voor dat de factor levels worden gerangschikt als "C1" (1), "C2"(2), "A" (3) en "B" (4).

Als u de rangorde van factor levels wilt omkeren, kan de vector van data values (gegevenswaarden) worden opgegeven als een argument voor de ingebouwde **rev()**-functie in de toewijzing aan het argument **levels=**:

- 1 Open RStudio en klik vervolgens op **File, New File, R Script** of druk op de toetsen **Ctrl + Shift + N** om een nieuw Code Editor venster te openen.
- 2 Voer de onderstaande broncode in de Code Editor in. Laat de groene commentaarregels weg.

```
FactorData.R* x
Source on Save
Run
Source

1 #Maak een data frame genaamed frame
2 frame <- data.frame( 1:5, sizes=c( "s", "L", "XL", "s", "M" ) )
3
4 #Voer de structuur van één kolom van het dataframe uit om de
5 #factoren en hun rangschikingsniveaus te zien
6 cat( "\nColumn Data...\n" )
7 str( frame$sizes )
8
9 #Maak een vector met dezelfde data values als de
10 #kolomgegevens, en in dezelfde volgorde
11 var.sizes <- c( "s", "L", "XL", "s", "M" )
12
13 #Voer de structuur van de vector uit om te zien dat deze alleen
14 #karaktergegevens bevat zonder factoren of levels
15 cat( "\nvector Data...\n" )
16 str( var.sizes )
17
18 #Reproduceer de vector om factoren en rangschikings-levels
19 #te creëren
20 var.sizes <- factor( var.sizes )
21
22 #Voer de structuur van de vector nogmaals uit om te zien dat
23 #deze nu factoren en levels heeft
24 cat( "\nFactored vector Data...\n" )
25 str( var.sizes )
26 print( levels( var.sizes ) )
27
28 #Reproduceer de vector om factoren en niveaus te creëren in een
29 #gewenste rangorde
30 var.sizes <- factor( var.sizes,
31                     levels = c( "s","M","L","XL" ) )
32
33 #Voer de herziene structuur van de vector levels uit
34 cat( "\nRe-ordered Factored Vector Data...\n" )
35 str( var.sizes )
36 print( levels( var.sizes ) )
37
38 #Reproduceer de vector om factoren en niveaus in een omgekeerde
39 #rangorde te maken en voer vervolgens de omgekeerde structuur uit
40 var.sizes <- factor( var.sizes,
41                     levels=rev( levels( var.sizes ) ) )
42 cat( "\nReversed Factored Vector Data...\n" )
43 str( var.sizes )
44 print( levels( var.sizes ) )
45
```



De functie **levels()** biedt toegang tot het attribuut levels van de variabele die als argument is opgegeven.

3

Sla het **R Script** op met de naam **FactorData.R** en klik vervolgens op de knop  **Source** of druk op **Ctrl + Shift + S** om de rekenkundige uitvoer te zien.

```
Console Terminal x Jobs x
C:/MyRScripts/
> source('E:/IT Concepten/IT Concepten - New Site/les
sen/rstudio/sourcecode/FactorData2.R')

Column Data...
chr [1:5] "S" "L" "XL" "S" "M"

Vector Data...
chr [1:5] "S" "L" "XL" "S" "M"

Factored Vector Data...
Factor w/ 4 levels "L","M","S","XL": 3 1 4 3 2
[1] "L" "M" "S" "XL"

Re-ordered Factored Vector Data...
Factor w/ 4 levels "S","M","L","XL": 1 3 4 1 2
[1] "S" "M" "L" "XL"

Reversed Factored Vector Data...
Factor w/ 4 levels "XL","L","M","S": 4 2 1 4 3
[1] "XL" "L" "M" "S"
> |
```



Factoren zijn unieke waarden. Dit voorbeeld heeft vijf gegevenswaarden, maar slechts vier unieke waarden - er