



Patrick Lürwer

# R4SEO: AUTOMATISIERTE REPORTS MIT R GENERIEREN (TEIL 3/4)

## DER AUTOR



**Patrick Lürwer** ist Senior Analyst bei get.traction GmbH. Dort ist er für die Datenerfassung, -aufbereitung und -analyse zuständig. Sein tägliches Handwerkszeug sind R, Python und KNIME.

Im vorherigen Teil dieser vierteiligen Artikelserie haben Sie gelernt, wie Sie fehlende Datumspunkte in Ihre DataFrames einfügen. Außerdem haben Sie die Berechnung des gleitenden Mittelwerts zur Glättung von Graphen angewendet und einen DataFrame transponiert. Damit bleiben noch zwei Aufgaben der Datenaufbereitung, bis Sie im nächsten Teil mit der Visualisierung beginnen können. Zum einen werden Sie in diesem Artikel die DataFrames mit den URLs resp. Queries aus der Google Search Console so aggregieren, dass Sie die Top-10-URLs bzw. -Queries nach Clicks in der Vorwoche erhalten. Zum anderen führen Sie die beiden DataFrames mit den Desktop- bzw. Mobile-Sichtbarkeitsdaten zusammen, um sie in einem Diagramm darstellen zu können.

```

189
190 # Daten aufräumen 02 -----
191
192
193 # GSC -----
194
195 # Berechnung des Montags und Sonntags der Vorwoche
196 MONDAY_PREVIOUS_WEEK ← floor_date(today(), "week", 1) - 7
197 SUNDAY_PREVIOUS_WEEK ← MONDAY_PREVIOUS_WEEK + 6
198
199
200 # Page:Berechnung der Top 10 Seiten nach Clicks der Vorwoche
201 gsa_top_10_pages ← gsa_dim_date_page %>%
202
203 1 filter(between(date, MONDAY_PREVIOUS_WEEK, SUNDAY_PREVIOUS_WEEK)) %>%
204
205   group_by(page) %>%
206   summarise(sum_clicks = sum(clicks),
207             2 sum_impressions = sum(impressions),
208             ctr = sum_clicks / sum_impressions,
209             weight_position = sum(impressions * position) / sum(impressions)) %>%
210   ungroup() %>%
211
212   3 top_n(n = 10,
213         wt = sum_clicks) %>%
214
215   4 arrange(desc(sum_clicks)) %>%
216
217   mutate(page = str_replace(page, GSC_PROP, "/"),
218          ctr = round(ctr, 3),
219          5 weight_position = round(weight_position, 2)) %>%
220
221   rename("Page" = page,
222          "Clicks" = sum_clicks,
223          "Impressions" = sum_impressions,
224          "CTR" = ctr,
225          "Position" = weight_position)
226

```

Abb.1: Berechnung der Top-10-Seiten nach Clicks der Vorwoche

```

1 # Hilfsfunktion
2 # print_weekday ← function(date) {
3   print(paste(date, "→", weekdays(date)))
4 # }
5
6
7 # Heutiges Datum
8 heute ← today()
9 print_weekday(heute)
10 # [1] "2019-01-20 → Sonntag"
11
12
13 # Datum des Wochenbeginns
14 beginn_der_woche ← floor_date(today(), "week")
15 print_weekday(beginn_der_woche)
16 # [1] "2019-01-20 → Sonntag" ← In den USA beginnt die Woche am Sonntag!
17
18
19 beginn_der_woche_montag ← floor_date(today(), "week", 1)
20 print(beginn_der_woche_montag)
21 # "2019-01-14 → Montag" ← 1 = Wochenbeginn am Montag
22
23
24 # Beginn der Vorwoche
25 beginn_der_vorwoche ← floor_date(today(), "week", 1) - 7
26 print_weekday(beginn_der_vorwoche)
27 # "2019-01-07 → Montag"
28
29
30 # Ende der Vorwoche
31 ende_der_vorwoche ← beginn_der_vorwoche + 6
32 print_weekday(ende_der_vorwoche)
33 # [1] "2019-01-13 → Sonntag"

```

Abb.2: Die Berechnung des Montags und Sonntags der Vorwoche im Detail

Im letzten Artikel haben Sie die Daten in den DataFrames *gsa\_dim\_date* sowie *ga\_sessions* transformiert. Damit fehlen noch die DataFrames für die URL- und Query-Performance-Daten (*gsa\_dim\_date\_page*, *gsa\_dim\_date\_query*) aus der Google Search Console (GSC) sowie den beiden DataFrames

für die Desktop- resp. Mobil-Sichtbarkeitsdaten (*si\_desktop*, *si\_mobile*) aus SISTRIX.

Zum Start legen Sie wie bisher ein neues R-Skript an, benennen es mit *data\_tidying\_02.R* und kopieren den Inhalt des Skripts *data\_tidying\_01.R* hinein. Sie müssen dieses Mal keine neuen

Packages installieren, sondern verwenden erneut jene, die Sie bereits aus dem vorherigen Artikel kennen. Führen Sie das neue Skript aus, damit wieder alle Daten zur Verfügung stehen.

## Top-Seiten aus der Google Search Console ermitteln

Für den zu erstellenden Report wollen Sie zwei Tabellen generieren, die die Top-10-Seiten bzw. -Suchanfragen nach Clicks der Vorwoche darstellen. Die dafür benötigten Daten liegen jedoch aktuell noch auf Tagesbasis für den gesamten Berichtszeitraum in den DataFrames *gsa\_dim\_date\_page* und *gsa\_dim\_date\_query*.

Kopieren Sie daher den Code aus Abbildung 1 ans Ende des zuvor von Ihnen neu erstellten Skripts. Keine Sorge, Sie werden diese Monsterabfrage gleich Schritt für Schritt durchgehen, um im Detail zu verstehen, wie die Daten transformiert werden. Die Absätze zwischen den einzelnen Abschnitten sind nur der Lesbarkeit geschuldet und können von Ihnen weggelassen werden.

Sie möchten für Ihren Report die Vorwoche betrachten. Daher müssen Sie zunächst ermitteln, auf welches Datum der Montag bzw. der Sonntag gefallen ist. Dazu dienen die ersten beiden Code-Zeilen. Anhand dieser können Sie dynamisch, ausgehend vom aktuellen Datum, die beiden Datumspunkte berechnen. Dabei gibt es eine Besonderheit: In den USA beginnt die Woche am Sonntag. Wenn Sie mit Zeit-

### TIPP

Um Ihnen das Abtippen des Codes zu ersparen, können Sie sich auch dieses Mal das Skript zu diesem Artikel unter folgendem Link auf GitHub herunterladen: [https://github.com/gettracti-ongmbh/r4seo\\_ws](https://github.com/gettracti-ongmbh/r4seo_ws). Ich empfehle Ihnen dennoch, den Code selbst zu schreiben, um ein Gefühl für Syntax und Funktionen zu bekommen.

reihen/-funktionen arbeiten, sollten Sie dies stets im Hinterkopf behalten. So wird bspw. auch in Google Analytics (GA) für die Diagrammdarstellung die amerikanische Woche genutzt. Wählen Sie im Kalender hingegen „Letzte Woche“, gilt wieder die deutsche Rechnung von Montag bis Sonntag.

Zurück zur Berechnung des Montags und Sonntags der Vorwoche. In Abbildung 2 sehen Sie die einzelnen Schritte im Detail. Die Funktion `floor_date()` aus dem Package `lubridate` kennen Sie bereits aus dem ersten Artikel. Dort haben Sie damit den Beginn des abzufragenden Berichtszeitraums berechnet. Mit dem Argument `week` wird der Funktion angezeigt, dass sie das Datum des Wochenbeginns, ausgehend vom heutigen Datum (`today()`), finden soll. Anschließend werden sieben Tage subtrahiert, um zum Montag der Vorwoche zu gelangen. Durch die Addition von sechs Tagen finden Sie dann den Sonntag. Die berechneten Datumsunkte schreiben Sie in die Variablen `MONDAY_PREVIOUS_WEEK` und `SUNDAY_PREVIOUS_WEEK`, um sie nachfolgend zum Filtern der DataFrames zu verwenden.

Im Anschluss führen Sie die Berechnung der Top-10-Seiten aus. Ihr DataFrame `gsa_dim_date_page` mit den Performance-Daten sieht wie in Abbildung 3 aus. Für jedes Datum des gesamten Berichtszeitraums, an dem die jeweilige URL in der Spalte `page` mindestens eine

Impression generiert hat, liegen Ihnen die Metriken `clicks`, `impressions`, `ctr` und `position` vor.

Zunächst müssen Sie den DataFrame auf den Datumsbereich der Vorwoche filtern (Abbildung 1; 1). Dazu geben Sie `gsa_dim_date_page` an die Funktion `filter()`. Diese nimmt wiederum die Funktion `between()` auf. In dieser definieren Sie, dass die Spalte `date` dahingehend geprüft werden soll, ob die Werte darin zwischen den beiden Datumsunkten `MONDAY_PREVIOUS_WEEK` und `SUNDAY_PREVIOUS_WEEK` liegen. Übrig bleiben somit nur noch die Zeilen, welche Daten der Vorwoche enthalten.

Darauf folgt ein `group_by()` (2), denn Sie wollen die Metriken für jede URL separat aggregieren. Wie die einzelnen Metriken aggregiert werden sollen, definieren Sie mit dem anschließenden `summarise()`. Die Operation ist etwas abstrakt, daher finden Sie in Abbildung 4 eine schematische Darstellung anhand von drei URLs. Der `group_by()`

teilt den DataFrame in einzelne Untergruppen basierend auf den URLs auf. Dies erfolgt nur intern, d. h., würden Sie den Code nur bis zum `group_by()` ausführen, sähen Sie keinen Unterschied in der Darstellung zum initialen DataFrame.

Innerhalb des `summarise()` definieren Sie dann, welche neuen Spalten gebildet werden sollen und wie die Werte dieser Spalten berechnet werden. Für die Spalten `clicks` und `impressions` berechnen Sie die Summe (`sum()`) jeder URL-Gruppe und schreiben das Ergebnis in die Spalten `sum_clicks` resp. `sum_impressions`. Die `ctr` berechnen Sie hier neu, indem Sie die Werte der neuen Spalte `sum_clicks` durch `sum_impressions` teilen. Für die `position` können Sie nicht einfach den Mittelwert berechnen, da die Wertigkeit einer Ranking-Position von den Impressions abhängig ist. Sie können auf Ihrer Website eine Seite haben, die zu einer sehr seltenen Suchanfrage auf Position 1 rankt. Hat diese Seite jedoch

| date       | page  | clicks | impressions | ctr        | position |
|------------|---|--------|-------------|------------|----------|
| 2018-09-19 | https://dein-trueffel.de/trueffel-produkte/trueffelbutter/    | 4302   | 38018       | 0.11315693 | 4.285286 |
| 2018-09-20 | https://dein-trueffel.de/trueffel-produkte/trueffelbutter/    | 787    | 6787        | 0.11595698 | 1.207603 |
| 2018-09-21 | https://dein-trueffel.de/trueffel-produkte/trueffelbutter/    | 393    | 3207        | 0.12254443 | 1.489866 |
| 2018-09-23 | https://dein-trueffel.de/trueffel-produkte/trueffelbutter/    | 263    | 1788        | 0.14709172 | 1.733781 |
| 2018-09-22 | https://dein-trueffel.de/trueffel-produkte/trueffelbutter/    | 246    | 2058        | 0.11953353 | 1.681244 |
| 2018-10-27 | https://dein-trueffel.de/trueffel-wissen/wo-wachsen-trueffel/ | 119    | 680         | 0.17500000 | 6.579412 |
| 2019-01-04 | https://dein-trueffel.de/trueffel-wissen/wo-wachsen-trueffel/ | 119    | 574         | 0.20731707 | 4.750871 |

Abb.3: DataFrame `gsa_dim_date_page`

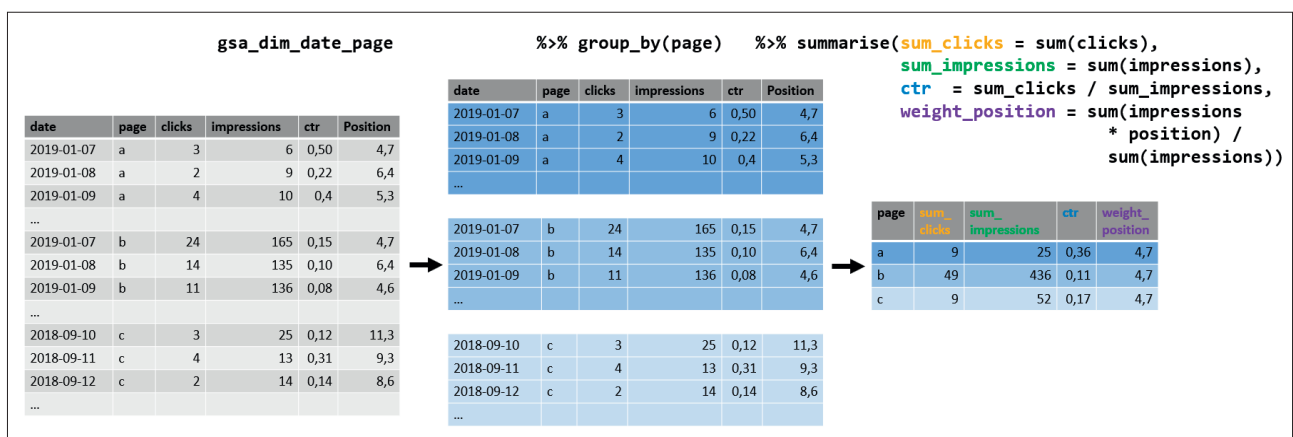


Abb.4: Schematische Darstellung von `group_by()` und `summarise()`

| page  | sum_clicks | sum_impressions | ctr        | weight_position |
|---|------------|-----------------|------------|-----------------|
| 1 https://dein-trueffel.de/   | 5          | 15              | 0.33333333 | 5.733333        |
| 2 https://dein-trueffel.de/trueffel-produkte/trueffelbutter/            | 107        | 1033            | 0.10358180 | 5.433688        |
| 3 https://dein-trueffel.de/trueffel-produkte/trueffelhonig/             | 4          | 6               | 0.66666667 | 5.666667        |
| 4 https://dein-trueffel.de/trueffel-produkte/trueffeloel/               | 2          | 20              | 0.10000000 | 18.150000       |
| 5 https://dein-trueffel.de/trueffel-rezepte/                            | 2          | 14              | 0.14285714 | 14.571429       |
| 6 https://dein-trueffel.de/trueffel-rezepte/383-3/                      | 4          | 44              | 0.09090909 | 36.909091       |
| 7 https://dein-trueffel.de/trueffel-rezepte/383-4/                      | 17         | 143             | 0.11888112 | 8.580420        |
| 8 https://dein-trueffel.de/trueffel-wissen/arten/china-trueffel/        | 14         | 87              | 0.16091954 | 9.643678        |
| 9 https://dein-trueffel.de/trueffel-wissen/arten/fruehlingstrueffel/    | 2          | 6               | 0.33333333 | 6.833333        |
| 10 https://dein-trueffel.de/trueffel-wissen/arten/perigord-trueffel/    | 10         | 243             | 0.04115226 | 9.226337        |
| 11 https://dein-trueffel.de/trueffel-wissen/arten/wintertrueffel/       | 6          | 79              | 0.07594937 | 10.822785       |
| 12 https://dein-trueffel.de/trueffel-wissen/trueffelsuche/              | 48         | 576             | 0.08333333 | 8.373264        |
| 13 https://dein-trueffel.de/trueffel-wissen/trueffelsuche/trueffelhu... | 4          | 125             | 0.03200000 | 13.056000       |
| 14 https://dein-trueffel.de/trueffel-wissen/trueffelsuche/trueffelsc... | 4          | 273             | 0.01465201 | 7.736264        |
| 15 https://dein-trueffel.de/trueffel-wissen/wo-wachsen-trueffel/        | 212        | 1178            | 0.17996604 | 6.238540        |

Abb.5: DataFrame nach der Aggregation auf URL-Ebene

| page  | sum_clicks | sum_impressions | ctr        | weight_position |
|---|------------|-----------------|------------|-----------------|
| 1 https://dein-trueffel.de/   | 5          | 15              | 0.33333333 | 5.733333        |
| 2 https://dein-trueffel.de/trueffel-produkte/trueffelbutter/            | 107        | 1033            | 0.10358180 | 5.433688        |
| 3 https://dein-trueffel.de/trueffel-produkte/trueffelhonig/             | 4          | 6               | 0.66666667 | 5.666667        |
| 4 https://dein-trueffel.de/trueffel-rezepte/383-3/                      | 4          | 44              | 0.09090909 | 36.909091       |
| 5 https://dein-trueffel.de/trueffel-rezepte/383-4/                      | 17         | 143             | 0.11888112 | 8.580420        |
| 6 https://dein-trueffel.de/trueffel-wissen/arten/china-trueffel/        | 14         | 87              | 0.16091954 | 9.643678        |
| 7 https://dein-trueffel.de/trueffel-wissen/arten/perigord-trueffel/     | 10         | 243             | 0.04115226 | 9.226337        |
| 8 https://dein-trueffel.de/trueffel-wissen/arten/wintertrueffel/        | 6          | 79              | 0.07594937 | 10.822785       |
| 9 https://dein-trueffel.de/trueffel-wissen/trueffelsuche/               | 48         | 576             | 0.08333333 | 8.373264        |
| 10 https://dein-trueffel.de/trueffel-wissen/trueffelsuche/trueffelhu... | 4          | 125             | 0.03200000 | 13.056000       |
| 11 https://dein-trueffel.de/trueffel-wissen/trueffelsuche/trueffelsc... | 4          | 273             | 0.01465201 | 7.736264        |
| 12 https://dein-trueffel.de/trueffel-wissen/wo-wachsen-trueffel/        | 212        | 1178            | 0.17996604 | 6.238540        |

Abb.6: Top-10-URLs nach Clicks

nur eine Impression generiert, sollte sie in Relation zu Ihren anderen Seiten mit mehr Impressions und ggf. schlechterem Ranking nicht als Ihre Top-Seite qualifiziert werden. Daher gewichten Sie für die Spalte *weight\_position* die durchschnittliche Position anhand der Impressions. Als Resultat erhalten Sie eine reduzierte Tabelle, in der je URL nur noch eine Zeile vorhanden ist. Die Datumsspalte entfällt automatisch.

Nach dem *summarise()* folgt noch ein *ungroup()*, welches die internen Gruppen wieder auflöst.

Der DataFrame sieht nun wie in Abbildung 5 aus. Für die Beispiel-Webseite sind insgesamt 15 URLs übrig geblieben, die in der vergangenen Woche Impressions generiert haben. Sie wollen in Ihrem Report aber nur die Top-10-URLs nach Clicks anzeigen. Zum Glück gibt es dafür eine *top\_n()*-Funktion (3). Das Argument *n* nimmt die Anzahl der zu behalten-

den Zeilen an. Mit dem Argument *wt* definieren Sie, anhand welcher Spalte die Auswahl der Zeilen erfolgen soll. Standardmäßig werden somit die zehn Zeilen mit den höchsten *sum\_clicks*-Werten ausgewählt. Wollen Sie die URLs mit den schlechtesten Click-Werten, können Sie ein Minus vor die Spalte schreiben (*wt = -sum\_clicks*). Weil einige URLs die gleiche Anzahl an Clicks erhalten haben, finden sich in der Tabelle allerdings zwölf statt zehn Zeilen.

Ihre Tabelle sieht zum jetzigen Zeitpunkt aus, wie in Abbildung 6 zu sehen. Da Sie die Tabelle in Ihrem Report gerne absteigend nach Clicks sortieren wollen, führen Sie im Anschluss die Funktion *arrange()* aus (4). In dieser geben Sie die zu sortierende Spalte an. Das *desc()* (descending) zeigt die absteigende Sortierung an. Lassen Sie es weg, erfolgt die Sortierung automatisch absteigend.

Nachdem Sie den DataFrame in Form gebracht haben, widmen Sie sich der Bearbeitung der Werte (5). Mit der Funktion *mutate()* können Sie bestehende Zeilen bearbeiten oder neue hinzufügen. Im vorliegenden Fall wollen Sie die Spalten *page*, *ctr* und *weight\_position* bearbeiten. Zuerst schneiden Sie die Domain aus den URLs ab. In der Tabelle des finalen Reports nimmt sie nur Platz weg. Sie bedienen sich dazu der Funktion *str\_replace()*. Als erstes Argument nimmt sie die Spalte auf, in der ein String ersetzt werden soll. Das zweite Argument *pattern* ist das Muster, das ersetzt werden soll. Hier geben Sie die Variable *GSC\_PROP* an, die Sie bereits im ersten Artikel dazu verwendet haben, die GSC-Daten von der API abzufragen; *replacement* nimmt dann das einzusetzende Zeichen an. Konkret ersetzen Sie dadurch bei allen URLs „*https://dein-trueffel.de/*“ durch „*/*“. Danach runden Sie die Werte der Spalten *ctr* und *weight\_position* auf drei bzw. zwei Nachkommastellen.

Zum Schluss bringen Sie noch etwas Hochglanz auf. Die Spaltennamen sind aktuell noch normalisiert, sprich, kleingeschrieben und mit Unterstrichen. Da das in Reports nicht unbedingt schön aussieht, benennen Sie sie mittels *rename()* um (6). Ihre finale Tabelle sehen Sie in Abbildung 7.

### Top-Queries aus der Google Search Console ermitteln

Geschafft! Nachdem Sie diese Tour de Force überstanden haben, lernen Sie einen der wesentlichen Vorzüge der Programmierung mit R kennen. Sie können Code-Abschnitte einfach wiederverwenden. Denn für die Berechnung der Top-10-Queries ist der Ablauf identisch (Abbildung 8). Sie müssen nur den anfänglichen DataFrame zu *gsa\_dim\_date\_query* ändern. Dann noch eine kleine Anpassung im *group\_by()*, denn jetzt aggregieren Sie ja auf Query- und nicht mehr auf

URL-Basis. Zu guter Letzt ändern Sie die Benennung des DataFrames, in den die Tabelle geschrieben werden soll, zu *gsa\_top\_10\_queries*. Die Tabelle mit den Top-Queries sehen Sie in Abbildung 9.

### S-DataFrames zusammenführen

Nun können Sie sich dem Zusammenführen der DataFrames mit den Sichtbarkeitswerten (SI) widmen. Fügen Sie dazu den Code aus Abbildung 10 in Ihr Skript ein. Der DataFrame *si\_desktop*, der die Sichtbarkeitsdaten der Desktop-Variante enthält, sieht aktuell aus, wie in Abbildung 11 zu sehen. Bevor Sie diesen DataFrame mit *si\_mobile* zusammenführen, gilt es noch drei Transformationen vorzunehmen (Abbildung 10; 1): Zunächst fällt auf, dass die *date*-Spalte das Datum mit Zeitstempel enthält. Diese sind für den späteren Report überflüssig und werden daher zu einem reinen Datumswert reduziert. Dazu greifen Sie erneut auf die Funktion *mutate()* zurück. In dieser wenden Sie die Funktion *date()*, die das Datum extrahiert, auf die Spalte *date* an. Den so gewonnenen Wert schreiben Sie in dieselbe Spalte zurück. Sie überschreiben somit die bestehenden Werte einfach. Anschließend entfernen Sie mittels *select()* die Spalte *domain*. Das vorangestellte Minus (-*domain*) zeigt eine negative Auswahl an, sprich, die Funktion soll alle Spalten außer *domain* auswählen. Schließlich benennen Sie durch *rename()* die Spalte *value* zu *Desktop* um. Dadurch machen Sie explizit, dass es sich bei den Werten dieser Spalte um die Desktop-SI handelt.

Würden Sie die Spalte nicht umbenennen und im nächsten Schritt die Mobile-Werte an die Tabelle schreiben, hätten Sie zwei Spalten *value* – und könnten folglich nicht ohne Weiteres sagen, welche Spalte die Desktop- resp. die Mobile-Werte enthält.

Und damit kommen Sie auch schon zur Funktion, die die beiden DataFra-

| Page   | Clicks | Impressions | CTR   | Position |
|--|--------|-------------|-------|----------|
| 1 /trueffel-wissen/wo-wachsen-trueffel/            | 212    | 1178        | 0.180 | 6.24     |
| 2 /trueffel-produkte/trueffelbutter/               | 107    | 1033        | 0.104 | 5.43     |
| 3 /trueffel-wissen/trueffelsuche/                  | 48     | 576         | 0.083 | 8.37     |
| 4 /trueffel-rezepte/383-4/                         | 17     | 143         | 0.119 | 8.58     |
| 5 /trueffel-wissen/arten/china-trueffel/           | 14     | 87          | 0.161 | 9.64     |
| 6 /trueffel-wissen/arten/perigord-trueffel/        | 10     | 243         | 0.041 | 9.23     |
| 7 /trueffel-wissen/arten/wintertrueffel/           | 6      | 79          | 0.076 | 10.82    |
| 8 /  | 5      | 15          | 0.333 | 5.73     |
| 9 /trueffel-produkte/trueffelhonig/                | 4      | 6           | 0.667 | 5.67     |
| 10 /trueffel-rezepte/383-3/                        | 4      | 44          | 0.091 | 36.91    |
| 11 /trueffel-wissen/trueffelsuche/trueffelhund/    | 4      | 125         | 0.032 | 13.06    |
| 12 /trueffel-wissen/trueffelsuche/trueffelschwein/ | 4      | 273         | 0.015 | 7.74     |

Abb.7: Die finale Tabelle *gsa\_top\_10\_pages*

```

230
231 # Query: Berechnung der Top 10 Queries nach Clicks der Vorwoche
232 gsa_top_10_queries <- gsa_dim_date_query %>%
233   filter(between(date, MONDAY_PREVIOUS_WEEK, SUNDAY_PREVIOUS_WEEK)) %>%
234   group_by(query) %>%
235   summarise(sum_clicks = sum(clicks),
236             sum_impressions = sum(impressions),
237             ctr = sum_clicks / sum_impressions,
238             weight_position = sum(impressions * position) / sum(impressions)) %>%
239   ungroup() %>%
240   top_n(10, sum_clicks) %>%
241   arrange(desc(sum_clicks)) %>%
242   mutate(ctr = round(ctr, 3),
243          weight_position = round(weight_position, 2)) %>%
244   rename("Query" = query,
245          "Clicks" = sum_clicks,
246          "Impressions" = sum_impressions,
247          "CTR" = ctr,
248          "Position" = weight_position)
    
```

Abb.8: Berechnung der Top-10-Queries nach Clicks der Vorwoche

| Query                      | Clicks | Impressions | CTR   | Position |
|----------------------------|--------|-------------|-------|----------|
| 1 trüffelbutter            | 62     | 563         | 0.110 | 2.75     |
| 2 wo wachsen trüffel       | 35     | 78          | 0.449 | 1.22     |
| 3 wo findet man trüffel    | 14     | 30          | 0.467 | 1.90     |
| 4 trüffelsuche             | 13     | 48          | 0.271 | 1.48     |
| 5 trüffel finden           | 8      | 21          | 0.381 | 2.38     |
| 6 wo wächst trüffel        | 8      | 15          | 0.533 | 1.40     |
| 7 trüffelbutter verwendung | 6      | 10          | 0.600 | 1.20     |
| 8 trüffel in deutschland   | 5      | 24          | 0.208 | 7.50     |
| 9 trüffel fundorte         | 4      | 8           | 0.500 | 1.12     |
| 10 trüffelbutter rezept    | 4      | 19          | 0.211 | 7.63     |

Abb.9: Die finale Tabelle *gsa\_top\_10\_queries*

mes miteinander verbindet (2). Konkret wollen Sie die Spalte mit den mobilen SI aus dem DataFrame *si\_mobile* als zusätzliche Spalte an den DataFrame *si\_desktop* schreiben, sodass eine Zeile je Datumswert den Desktop- und den Mobile-SI listet. Vergleichbar ist diese Operation, die man als Join bezeichnet, mit einem *SVERWEIS* in Excel.

Sie ist allerdings etwas komplexer, denn die beiden DataFrames weisen unterschiedliche Datumsbereiche auf. Der Desktop-SI wird seit 2008, der Mobile-SI erst seit 2015 reportet. Für jeden Datumspunkt des Desktop-SI gibt es somit nicht unbedingt eine entsprechende Zeile im Mobile-SI. Da Sie im Rahmen von Analysen sehr häufig

```

250
251 # Sistrix -----
252
253 si <- si_desktop %>%
1  mutate(date = date(date)) %>%
   select(-domain) %>%
254   rename("Desktop" = value) %>%
255
256
257
258 left_join(si_mobile %>%
2  mutate(date = date(date)) %>%
   select(-domain) %>%
   rename("Mobile" = value),
262   by = "date") %>%
263
264
265
266
267
268
269
3  replace_na(list(Mobile = 0)) %>%
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

Abb.10: Zusammenführen der Sichtbarkeits-DataFrames

Tabellen miteinander verbinden werden, finden Sie in Abbildung 12 eine Übersicht der gängigsten Joins.

Im vorliegenden Beispiel wissen Sie, dass der DataFrame *si\_desktop* mehr Zeilen enthält als *si\_mobile*. Sie wenden daher einen *left\_join()* an, denn

Datumspunkte in *si\_desktop*, die keine Entsprechung *si\_mobile* aufweisen, sollen nicht entfallen (dies wäre bei einem *inner\_join()* der Fall), sondern zunächst mit NAs aufgefüllt werden und damit erhalten bleiben. Der *left\_join()* ist nun wie folgt aufgebaut: Zunächst nimmt

### HINWEIS

Haben Sie Fragen? Dann stellen Sie sie gerne in der Facebook-Gruppe <https://www.facebook.com/groups/ompyr/>. So können alle von den Antworten profitieren. Seien Sie versichert, es gibt keine „dummen“ Fragen, denn aller Anfang ist schwer.

die Funktion als erstes Argument den DataFrame auf, der mit der *si\_desktop* zusammengeführt werden soll. Dies ist der *si\_mobile*. Verschachtelt, innerhalb des *inner\_join()*, transformieren Sie diese Tabelle analog zu *si\_desktop*, mit dem Unterschied, dass Sie die Spalte *value* nun zu *Mobile* umbenennen. Das Argument *by* gibt an, welche Spalte als Schlüssel (Key) dienen soll, anhand dessen die beiden DataFrames miteinander gejoint werden. Hier wählen Sie die Spalte *date*, denn Sie wollen,



Die erschreckenden Risiken unserer vernetzten Welt und die beängstigende Wahrheit über den Sicherheitszustand des »Internet of Things«

384 Seiten | 24,99 € (D)  
ISBN 978-3-95845-947-2  
[www.mitp.de/947](http://www.mitp.de/947)



Mit realen Fallbeispielen von der Profilerstellung zu Werbezwecken bis hin zur Überwachung – und wie Sie sich dagegen schützen

320 Seiten | 24,99 €  
ISBN 978-3-95845-635-8  
[www.mitp.de/635](http://www.mitp.de/635)



Auswirkungen Künstlicher Intelligenz auf unser Leben, die Wirtschaft und soziale Strukturen sowie die einfache Einführung in die Grundlagen

208 Seiten | 24,99 €  
ISBN 978-3-95845-632-7  
[www.mitp.de/632](http://www.mitp.de/632)

dass für jedes Datum der Desktop- und der Mobile-SI in einer Zeile stehen. Die zunächst mit *NA* aufgefüllten Werte in der Mobile-Spalte ersetzen Sie anschließend durch die bereits bekannte Funktion *replace\_na()* mit 0 (Abbildung 10; 3).

Für Ihren Report haben Sie definiert, dass Sie nur einen Berichtszeitraum von rückwirkend sechs Monaten verwenden möchten (*START\_DATE* und *END\_DATE* aus Teil 1 stehen am Anfang des Skripts). Über die Filterfunktion schränken Sie den DataFrame auf diesen Bereich ein (4). Damit machen Sie im vorliegenden Beispiel den vorherigen Schritt des Auffüllens der mobilen SI-Werte natürlich im Grunde überflüssig, da Sie die Zeilen direkt wieder verwerfen. Dass Sie diesen Schritt jedoch durchgeführt haben, bietet den Vorteil, dass Sie fortan den Berichtszeitraum nach Belieben in die Vergangenheit vergrößern können.

Zum Schluss überführen Sie den nun weiten DataFrame mittels *gather()* in einen langen (5), damit er eine zum Plotting geeignete Struktur hat. Der finale DataFrame *si* sieht nun aus, wie in Abbildung 13 zu sehen.

### Fazit

Jetzt bleibt nur noch eines übrig: Klopfen Sie sich kräftig auf die Schulter und vergegenwärtigen Sie sich noch einmal, was Sie alles geschafft und gelernt haben!

Sie haben die Top-10-Seiten nach GSC-Clicks berechnet. Dabei haben Sie die zwei zentralen Funktionen zur Datenaggregation kennengelernt: *group\_by()* und *summarise()*. Kein Report und keine Analyse kommen ohne diese beiden aus, denn sie ermöglichen es Ihnen, einen DataFrame entsprechend den Ausprägungen in einer Spalte zu gruppieren und die einzelnen Gruppen gesondert voneinander zu aggregieren. Ein simples Beispiel aus einem anderen Bereich ist das Zählen

|   | domain           | date                      | value  |
|---|------------------|---------------------------|--------|
| 1 | dein-trueffel.de | 2019-01-28T00:00:00+01:00 | 0.0002 |
| 2 | dein-trueffel.de | 2019-01-21T00:00:00+01:00 | 0.0002 |
| 3 | dein-trueffel.de | 2019-01-14T00:00:00+01:00 | 0.0015 |
| 4 | dein-trueffel.de | 2019-01-07T00:00:00+01:00 | 0.0002 |
| 5 | dein-trueffel.de | 2018-12-31T00:00:00+01:00 | 0.0002 |
| 6 | dein-trueffel.de | 2018-12-24T00:00:00+01:00 | 0.0002 |
| 7 | dein-trueffel.de | 2018-12-17T00:00:00+01:00 | 0.0000 |

Abb.11: DataFrame si\_desktop

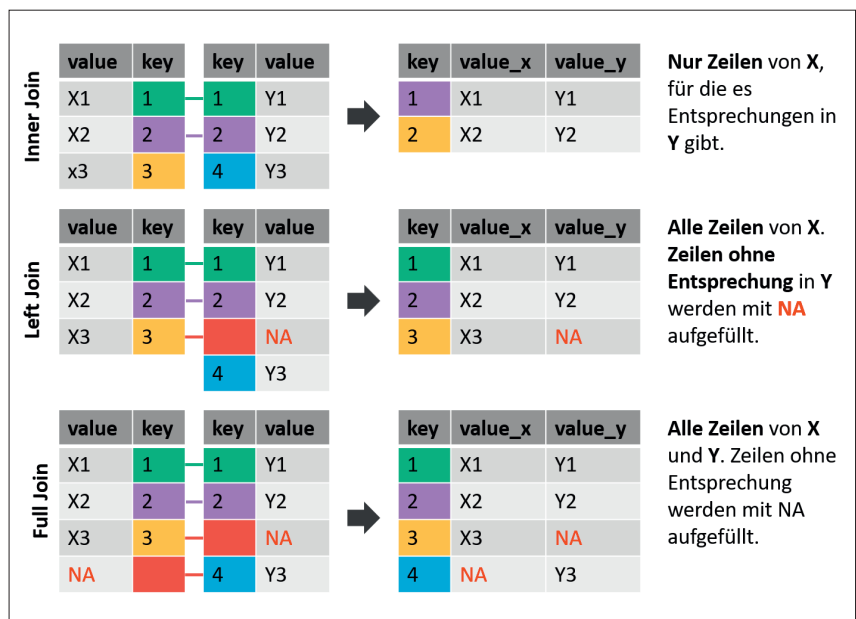


Abb.12: Übersicht von Join-Operationen

der einzelnen Status-Codes in einem Crawl oder die Berechnung der durchschnittlichen Verweildauer auf den Seiten der einzelnen Kategorien in einem Blog.

Sie haben einen kurzen Abstecher in die String-Manipulation absolviert, indem Sie mittels *str\_replace()* die Domain aus den URLs entfernt haben. Replace ist eine mächtige, aber bei Weitem nicht die einzige Funktion in diesem Bereich. Geben Sie einmal *str\_* auf der Console ein und scrollen Sie durch die Autovervollständigung. Sie finden eine Fülle von Funktionen, um

Bestandteile von Strings zu entfernen, zu transformieren oder zu extrahieren. Gerade beim Arbeiten mit URLs ist das ein unverzichtbares Werkzeug, um bspw. irrelevante Parameter abzuschneiden, URLs der mobilen auf die Desktop-Variante umzuschreiben oder Produkt-IDs zu extrahieren.

Schließlich haben Sie einen Join zweier DataFrames durchgeführt – ebenfalls ein elementares Werkzeug bei der Arbeit mit Daten, denn fast immer müssen diese aus verschiedenen Quellen zusammengeführt werden. ¶