

Testat Webtext

Patrick Bucher

06.04.2017

Inhaltsverzeichnis

Linux: Produktiver arbeiten mit Textdateien	1
Anwendungsbeispiel: Wörter in mehreren Artikeln zählen	1
Der wc-Befehl	2
Die Ausgabe an sort weiterleiten	2
Die Ausgabe kürzen	3

Linux: Produktiver arbeiten mit Textdateien

Linux-User mögen Textdateien. Gegenüber binären oder XML-basierten Formaten wie *OpenDocument Text* oder *Office Open XML* haben sie einige Vorteile:

- Sie lassen sich mit einem beliebigen Editor bearbeiten.
- Sie benötigen viel weniger Speicherplatz.
- Sie können auch in ferner Zukunft noch gelesen werden.
- Sie lassen sich als Textstrom bearbeiten.

Doch was hat man als Benutzer davon, wenn man eine Datei als *Textstrom* bearbeiten kann?

Anwendungsbeispiel: Wörter in mehreren Artikeln zählen

Angenommen, wir haben eine Reihe von Artikeln; einmal im *OpenDocument Text*-Format (.odt) und einmal im Textformat (.txt). Nun wollen wir herausfinden, welcher Artikel in Wörtern gemessen der längste ist. Mit unseren *OpenDocument Text*-Dateien verfahren wir folgendermassen:

1. Wir öffnen die erste Datei mit *LibreOffice* oder *OpenOffice.org*.
2. Wir gehen auf das Menü "Tools" und wählen den Eintrag "Word Count".
3. Wir notieren uns den Dateinamen und die Anzahl Wörter dazu.
4. Wir schliessen die Datei und fahren für die nächste Datei bei Schritt 1. fort.

Dieses Vorgehen ist sehr aufwändig. Zudem muss der ganze Vorgang zu einem späteren Zeitpunkt wiederholt werden, falls die Dateien in Zwischenzeit bearbeitet wurden. Schliesslich könnte sich dadurch die Anzahl der Wörter verändert haben.

Mit Textdateien funktioniert das einfacher. Man verwendet einfach folgenden Konsolenbefehl bzw. folgende zwei, durch eine sogenannte *Pipe* verbundenen Befehle:

```
$ wc -w *.txt | sort -n -r
```

Dadurch erhält man folgende Ausgabe:

```
2220 eigenes-bier-brauen.txt
1739 berlinreise.txt
1231 neues-aquarium.txt
893 im-stau.txt
```

Doch was hat das ganze zu bedeuten? Schauen wir uns den Befehl (bzw. die Befehle) genauer an:

Der `wc`-Befehl

- Das `wc` zählt die Wörter in einer Datei. Der Programmname ist eine Abkürzung für “word count”.
- Standardmässig gibt `wc` die Anzahl Zeilen, Wörter und Zeichen einer Datei aus. Wir interessieren uns aber nur für die Wörter. Darum geben wir den Parameter `-w` (für “words”) mit.
- Mit `*.txt` geben wir dem Programm sämtliche Textdateien im aktuellen Arbeitsverzeichnis zum Zählen.

Das (`wc -w *.txt`) ist der erste Teil des Befehls. Führt man ihn aus, erhielte man folgende Ausgabe:

```
1739 berlinreise.txt
2220 eigenes-bier-brauen.txt
893 im-stau.txt
1231 neues-aquarium.txt
```

Die Ausgabe an `sort` weiterleiten

Die Dateien sind alphabetisch und nicht nach Nummern sortiert. Das liegt daran, dass `wc` die Dateien über die Wildcard `*.txt` in alphabetischer Reihenfolge zur Bearbeitung erhält. Darum kommt jetzt der zweite Teil der Befehlszeile zum Zug:

- Zwischen den Befehlen steht `|`: eine sogenannte *Pipe*, zu Deutsch etwa “Röhre”. Diese nimmt die Ausgabe eines Programmes entgegen und leitet sie als Eingabe zum nächsten Programm weiter.
- Das nächste Programm ist in diesem Fall `sort`, das Textzeilen in alphabetisch und in aufsteigender Reihenfolge sortiert.
- Da wir aber keine alphabetische, sondern numerische Sortierung wollen (“100” wäre gemäss alphabetischer Sortierung kleiner als “9”), geben wir den Parameter `-n` an.

- Zudem soll die Reihenfolge nicht aufsteigend (die kleinste Zahl am Anfang) sondern absteigend (die grösste Zahl am Anfang) sein, was wir mit dem Parameter `-r` machen.

Die Ausgabe kürzen

Je mehr Artikel sich in unserem Verzeichnis befinden, desto länger wird die Ausgabe. Bei hunderten Artikel müssten wir bald nach oben scrollen, um zu sehen, welcher am meisten Wörter enthält. Nun könnte man natürlich die Sortierreihenfolge anpassen, sodass der Artikel mit den meisten Wörtern in der letzten Zeile steht. Dazu kann man einfach den Parameter `-r` beim `sort`-Befehl weglassen.

Eleganter ist es, die Ausgabe auf eine bestimmte Länge zu kürzen. Der `head`-Befehl nimmt beliebig viele Zeilen entgegen und gibt die per Parameter definierte Anzahl an Zeilen aus. Um die drei Artikel mit den meisten Wörtern zu ermitteln, kann man den Befehl `head -3` verwenden. Dieser Befehl wird wiederum über eine Pipe mit unserer Befehlszeile verbunden:

```
$ wc -w *.txt | sort -n -r | head -3
```

Die Ausgabe könnte dann folgendermassen aussehen:

```
9431 der-syrienkrieg.txt
8943 jahresrückblick-2016.txt
7131 tour-de-france_doping.txt
```