

The Ninth International Conference on Advances in Databases, Knowledge, and Data Applications, Mai 21 - 25, 2017 - Barcelona, Spain

Tutorial: Data Manipulation and Data Transformation using the Shell  
 Andreas Schmidt, Steffen Scholz

## Command Overview

| Command | Description  | Popular options  |
|---------|--|--|
| cat     | concatenate files and print on the standard output   | -n: number all output lines  |
| head    | output the first part of files   | -n<num>: print the first <num> lines<br>-n -<num>: print all but the last <num> lines  |
| tail    |  | -n<num>: print the last <num> lines<br>-n + <num>: print, starting from line <num>   |
| wc      | print newline, word, and byte counts for each file   | -c: print byte counts<br>-m: print the character counts<br>-w: print the word counts<br>-l: print the newline counts   |
| grep    | print lines matching a pattern   | -E: support extended regexp<br>-i: ignore case<br>-v: invert match<br>-o: print only the matched part (output: one line per match)<br>-f: obtain patterns from file<br>-l: files with match<br>-L: files without match<br>-c: suppress normal output, instead count matching lines.<br>-m<num>: Stop reading a file after <num> matching lines |
| seq     | print a sequence of numbers  | -s: separator (default: \n)  |
| split   | split a file into pieces   |  |
| cut     | remove sections from each line of files  | -d<delim>: Use <delim> instead of ab as field separator<br>-f<field-list>: select only fields in <field-list><br>--output-delimiter=<delim> : use <delim> as output delimiter  |
| paste   | merge lines of files   | -d<char>: use <char> as output delimiter   |
| tr      | Translate, squeeze, and/or delete characters   | -c: use the complement of set1<br>-d: delete the characters in set1<br>-s: replace each sequence of a repeated character with a single character   |
| sort    | sort lines of text files   | -n: numeric sort<br>-r: reverse sort<br>-R: random shuffle<br>-c: check, if sorted, do not sort<br>-t: field separator<br>-k<keydef> : sort according to keydef<br><keydef>: F[.C][OPTS],[F[.C][OPTS]]<br>-u: output only the first of equal lines   |
| join    | join lines of two files on a common field<br><br>Remark: files must be sorted on join column | -t<char> : Use <char> as input, output separator<br>-1<field> : join on this FIELD of file 1<br>-2<field> : join on this FIELD of file 2<br>-o<format> : obey <format> while constructing output line  |

| Command | Description                           | Popular options   |
|---------|---------------------------------------|---|
| comm    | compare two sorted files line by line | <b>-1:</b> suppress column 1 (lines unique to FILE1)<br><b>-2:</b> suppress column 1 (lines unique to FILE2)<br><b>-3:</b> suppress column 1 (lines unique to FILE3)<br><b>--total</b> : output a summary |
| uniq    | report or omit repeated lines         | <b>-c</b> : prefix lines by the number of occurrences<br><b>-d</b> : only print duplicate lines, one for each group<br><b>-i</b> : ignore case<br><b>-u</b> : only print unique lines                     |

### sed – stream editor for filtering and transferring text

| Command | Description                                       | Popular options   |
|---------|---|---|
| sed     | stream editor for filtering and transforming text | <b>-n</b> : suppress automatic printing of pattern space<br><b>-f &lt;script-file&gt;</b> : scripts with commands to be executed<br><b>-i</b> : edit inplace<br><b>-E, -r</b> : support extended regexp   |
|         |   | <address><br><start-address>,<end-address><br><start-address>, + <number-of-lines><br><br><address> can be: <ul style="list-style-type: none"> <li>• line-number (i.e. 1,5,7, ...)</li> <li>• \$ (represent last line of file)</li> <li>• regular-expression</li> </ul> |

| Sed commands          | Description                          |
|-----------------------|--------------------------------------|
| a <text>              | append text                          |
| i <text>              | insert text                          |
| c <text>              | replace the selected lines with text |
| p                     | print                                |
| d                     | delete pattern space                 |
| s/regexp/replacement/ | regexp-replace                       |

### sed-Examples:

- sed -i '/Aachen/ d' city.csv # delete line(s) containing Aachen (inplace)
- sed '2i Karlsruhe,D,"Baden Wuerttemberg",301452,49.0,6.8' city.csv # insert 'Karlsruhe ...' at line 2
- sed -Ei '/<script>/,/</script>/d' jaccard.html # remove all script-sections
- sed -i 's/\bNULL\b/\N/g' city.csv # replace NULL ->\n
- sed -n '5,10p;23p;56,71p' city.csv # print lines 5-10, 23, 56-71