

Playing with Text and Data Files

Robbe Wüschendorf
10.1007/978-3-642-048-8-7
© Springer International

Abstract

How to look at gigabyte large text files? How to extract relevant or unique lines? How to compare text files? And, how to edit them? This chapter introduces basic command line tools and introduces you to the powerful Vim text editor.

7.1 Introduction

After learning some basics in the previous chapter, you will learn more advanced tools for editing text files and how you can use them to process large text files. In this chapter, that is, a text, certain program will not run immediately. This means you need to edit text files, too. In this chapter you will use very basic text editing tools. We will concentrate on the text editor Vim.

With Linux you have the choice from an endless list of text editors. From Microsoft Windows you might have heard about the Notepad. That is the standard Windows text editor. With Linux you can choose between `ed`, `vi`, `vim`, `emacs`, `sed`, `head`, `tail`, `cat`, and many more. In general, we can distinguish between three different types of text editors: (1) One group of text editors is called line-oriented. Vim belongs to this group. You can only work on one line of the text file and then need a command to get to the next line. This means that you cannot use the arrow keys to move the cursor (your current position) through the text. This is a little bit tedious and good only for editing pure text files (e.g., text files of editor source code). Vim belongs to this group of editors. You use the text over the whole screen and can scroll up and down and move to every text position using the arrow keys. This is much more comfortable than using a line-oriented editor. We will predominantly work with Vim. It is convenient and controllable enough for our purposes and usually available on all systems. You can even use it when you have no X-Server (see Sect. 3.3 on p. 26) running, which is often the case when you login to a remote computer. (2) The most comfortable text editors, of course, use the X-Server. However, you should not confuse comfortable with powerful. `sed` or `awk` are examples of editors which use the graphical user interface (GUI). Here, you can use the mouse to jump around in the text and copy and paste (drag and drop) marked text. We are going to work with these two editors.

7.2 Viewing and Analyzing Files

This section introduces some important tools that allow you to analyze text-based files in the command line.

7.2.1 A Quick Start: `cat`

Now, let us start with a very easy method to input text. For this, we use the command `cat` (concatenate). We have already used `cat` for displaying the content of a text file (e.g., for example Terminal 1 on p. 46). We have also learned about redirecting the output of `cat` into a file by using `>`. Now, when you use `cat` without specifying a filename you redirect the standard input (the keyboard) to the standard output (the screen). Let us try it out.

```
Terminal 24: Printing text with cat
1 cat
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

In Terminal 24 we first run the command `cat` without anything. After hitting `Enter` we get into line 2 and enter the text "this is text, now press Enter". At the end of line 2 we press `Enter`. The text we just entered will be printed again and we end up in line 4. Here we enter "now press Ctrl-C" and then press `Ctrl` + `C`. Again, the text will be printed, but now we end up in the same line. Now we press `Ctrl` + `Enter` and then `Ctrl` + `D`. We are back. What happens is that the input from the keyboard is redirected to the screen. We use `cat` to print `Enter` a new line from the keyboard. But this is not important now. It is important that you realize that we could also redirect the standard input into a file. How?

```
Terminal 25: Printing text with cat
1 cat
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

In Terminal 25 we first run the command `cat` without anything. After hitting `Enter` we enter some text and then hit `Ctrl` + `Enter`. We end up in line 5, enter some text again and press `Ctrl` + `D`. We have back the command line directly. No new line is generated. This is uncomfortable! If we now, as shown in lines 3–5, display the text with `cat test.txt`, we get the same phenomenon in line 5: the dollar character `$` is on the line. We have learned a new way of redirecting a file to a program. You now press `Enter` to get into a new file. Here, you should always finish a text file with a new line. Finally, in line 6 we remove the file that we have just created with the command `rm test.txt`.

It is easy to create a text file, is it not? Now create a file named `genomes.txt` as shown in Table 1.

```
Terminal 26: Creating text file
1 cat > genomes.txt
2 this is text, now press Enter
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Now let us see what we can do with the text file. We have already learned something about sorting the content of a text file in Terminal 9 on p. 48. Let us take a look at it again.

7.2.2 Text Sorting

With `sort` you can, as the name implies, sort the content of a text file. This can be very helpful in order to make files more readable. The `sort` command comes with some useful options, which are, as usual, separated by a dash.

- `Sorts` lines into account the value of numbers. Otherwise, only their first field (the alphabet would be considered and `D` would be printed before `E`).
- The `Sort` result is displayed in reverse order.
- Lower and uppercase characters are treated equally.
- Identical lines are printed only once. This you get unique or non-redundant output (see Sect. 7.1.3 on p. 61).
- Sorts according to the content of field `x`. For example, `-k 3` sorts according to the content of field 3.
- Set as a new field separator (instead of the space character `' '`). For example, `-t ','` sets the colon and `-t '\t'` sets the tabulator as field separator, respectively.

How does this work in real life?

```
Terminal 27: Sorting text with sort
1 sort -n genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

In Terminal 26 we use the options `-n` and `-k 3` in order to sort the content of the `genomes.txt` file, take care of correcting sorts of numbers and get rid of doublets, but what is this? The output of `sort -n genomes.txt` is displayed in reverse order. Three genomes are missing. Okay, here we encounter an error (bug) in the program `sort`. With such a well-established command as `sort`, this is a rather rare situation. What can we do? First of all, you should report the bug to the Linux community and then help to improve the command. In the manual pages (see Sect. 2.4.6 in slightly abbreviated form) you find the command `report-bug`. There is a statement where it is shown how to report bugs: "Report bugs to bug.debian@gnu.org". Then you have to think about alternatives to your original task. In our case the alternative is to sort in line 5. Here, we learn a new way of redirecting a file to a program. You now press `Enter` to get into a new file. Here, you should always finish a text file with a new line. Finally, in line 6 we remove the file that we have just created with the command `rm test.txt`.

The next example uses the options `-t ','` and `-k 1` in order to define a new field delimiter and the field, which is to be considered for sorting.

```
Terminal 27: Field-oriented sorting with sort
1 sort -t ',' -k 1 genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

With the `sort` command in line 1 of Terminal 27 the file `genomes.txt` is sorted according to the contents of field 3 (`-k 3`). The field delimiter is set to the dash character (`-t ,`). With these options, `sort` becomes very versatile and can be used to sort tables according to the values of a certain column.

7.2.3 Extract Unique Lines

Frequently, it is desired to extract unique lines from a text file. One example would be that you delete duplicate gene IDs from a file. The command `sort` can do with the option `-u`. Let us wait to show an alternative to this with the `uniq` command. Terminal 28 shows an example.

```
Terminal 28: Extracting unique lines with sort and uniq
1 sort genomes.txt | uniq
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

The `uniq` command requires that the data lines are sorted. I will replace successive lines containing a character different with just one copy of this line.

7.2.4 Viewing File Beginning or End

Imagine you have a number of text files but you cannot remember their content. You use `ls` to list the files. You could use `ls -l` to see the permissions and the size of the files. Or, of course, you can use `cat` to view the content of a file. But if you have a large file, you might not want to see the whole file and scroll through it if you want to read a part of the file on the screen, you should use `less`. (There is another program called `more`. It is older and thus not powerful. Just to let you know.) The syntax is easy:

```
Terminal 29: Viewing file beginning and end
1 less genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

After invoking `less`, the next page is displayed with `Enter` and the next line is displayed with `Enter`. By pressing `q`, you quit `less`. After typing the slash character `/`, you can enter a query term. With `h` and `j`, you jump to the next and previous lines matching the query pattern, which can be a regular expression (see 11 on p. 157). Another nice function is to jump to the end of the file with `g` or `G` to the beginning of the file with `1` or `G`. Accordingly, you can jump to any line, e.g., the 34 by typing `34` or `G`.

On some systems you will find the command `less`, which can be used in order to view compressed files.

7.2.6 Character, Word, and Line Counting

Another thing we might want to do is count the number of lines, words, and characters of a text file. This task can be performed with the command `wc` (word count). `wc` counts the number of bytes, whitespace-separated words and lines in each given file.

```
Terminal 30: Counting Lines and Words and Characters with wc
1 wc genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

As we can see from the output in Terminal 30, the file `genomes.txt` consists of 7 lines, 44 words, and 247 characters.

7.2.7 Splitting Files into Pieces

The `split` command splits files into a number of smaller files. If the list of output files is too long, you can use the option `-l` to limit the number of lines in each file. Terminal 31 shows an example.

```
Terminal 31: Splitting files with split
1 split genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

The `split` command in Terminal 30 on the preceding page split the file `genomes.txt` into small files. Each resulting file contains just lines 1 to 7 and the prefix of the resulting files is `xx`. The appendices `aa`, `ab` and so forth are also possible.

7.2.8 Cut and Paste Columns

You can most probably use `cut` and `paste` when writing with office applications or the like. This is where it is different. The `cut` cuts columns of a file. Columns are defined as text strings that are separated by a delimiter defined with `-d`. The option `-f` (fields) states the columns that shall be extracted. Terminal 32 shows how to work.

```
Terminal 32: Cutting columns with cut
1 cut -d ',' -f 1 genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

The command `paste` on the other hand fuses files to a table contained in one single file. Terminal 33 shows how.

```
Terminal 33: Pasting files to columns with paste
1 paste genomes.txt genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Assuming that you have three files with two lines of content each `paste` would fuse the file contents to a single tab-delimited file. The option `-s` makes `paste` reading the files in serial order.

7.2.9 Finding Text: `grep`

A very important command is `grep` (global regular expression printer), for a historical background see Sect. 11.1 on p. 157. `grep` searches the named input file(s) for lines containing a match to a given pattern. By default, `grep` prints the matching lines. Terminal 34 illustrates its function.

```
Terminal 34: Finding lines with grep
1 grep -n "AT" genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

The name `grep` has nothing to do with the verb to grab. It is derived from a function of a very old Unix text editor called `ed`. Terminal 34 shows how the text `AT` can be operated with `ed` (line 4). The command line now accepts `ed` commands (I added contents of the `ed` manual in the `ed` manual – you should not rely on it).

```
Terminal 35: Finding lines with ed
1 ed genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

The result of the command `ed` (different) indicates what lines have to be with the file `amorf` to convert it into amino. You have to delete the lines marked with `c` and add the lines marked with `C`. The option `-n` (line 18) the same text is shown in another way. The option `-c` (line 28) is used to display the differences only (indicated by `+`).

The command `--sort=command` you give the input files are sorted. We do this in line 1 in Terminal 36. Note: we write several commands in one line, separating them with the semicolon character (`;`).

```
Terminal 36: Finding things in command with sort
1 sort -n genomes.txt | grep -n "AT"
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

The `comm` command prints out its result in three columns. Column one contains all lines that appear only in the first file (`amorf`), column two shows all lines that are present only in the second file (`amorf`), and the third column contains all lines that are in both files. You can restrict the output with the option `-2`, which shows only the second column (minus 1 and minus 2), which contains the common file content.

7.3 Editing Text Files

It is very helpful to know how to edit files in the command line. Suppose you are connected to a remote computer without graphical interface. If you have the file `genomes.txt` in your home directory, you can use `cat` to view the content of the file. But if you have a large file, you might not want to see the whole file and scroll through it if you want to read a part of the file on the screen, you should use `less`. (There is another program called `more`. It is older and thus not powerful. Just to let you know.) The syntax is easy:

```
Terminal 37: Editing text with vi
1 vi genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Vi (pronounced vee-eye improved, or just vi) provides text editing capabilities for beginners and experts alike. Three aspects of Vi make it a powerful text editor. First, it is supplied with all Unix systems. Thus, you can use Vi on all universities or on any business with Linux or Unix systems. Second, Vi uses a small amount of memory, which allows efficient editing even when the network is busy. Third, Vi uses a simple standard alphanumeric key for commands, you can use it on virtually any terminal or workstation in existence without having to worry about unusual keyboard mappings.

You see, there are really many advantages to learning Vi. In celebration of the twentieth birthday of Vi Vi has even been ported to Apple's OS (see <http://www.apple.com/applelink/applelink.cfm?applelink=applelink.cfm>) and other mobile devices, e.g., Android (<http://play.google.com/store/apps/details?id=com.espyrid.vim&hl=de>). Additionally, it has many features of the native version such as automatic indentation, a visual mode, language-aware syntax highlighting, integrated spelling and grammar checking, and so on.

In case that you are working with Ubuntu, e.g., as outlined in Sect. 4.1.2 on p. 49, you need to install Vi. By default, Ubuntu comes with Vi installed. However, Vi is much more comfortable to use. Therefore, execute the following command in the terminal to install Vi:

```
sudo apt-get install vim
```

7.3.4 Immediate Takeoff

If you cannot wait to use Vi and make your first text file: here you go. Type `vi` to start the editor. Press `h` to start the editor mode and enter your file. You start in line 1 with `h`. You can save and quit at once with `:w` or `:q`. If you do not want to save your changes you must force quitting with `!q`. Be cautious when abandoning Vi in this manner! Because any changes you have made will be permanently lost.

Up to this point you have learned more than enough commands to use Vi in a powerful and controllable enough for our purposes and usually available on all systems. You will wish to use:

```
Terminal 38: Editing text in Vi
1 vi genomes.txt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Commands are very often preceded by the colon `:` character. Let us try another command: type `set number`. Now you see line numbers in front of each line. With `set number` you hide them again. Another command: type `if !a`. With `if !a`, you have a list of all files in your current working directory imported to your text file. This is magic. Is it not? Vi can search the entire file for a pattern and import the result to the text file at the current cursor position.

Probably `if !a` is not what you want. Let us try `if !a`. Usually you can use the arrow keys (if you use the cursor keys and letters instead of a moving cursor) to do the following: `h` twice; type a column `b`; type `wc` (non-capitalized); press `Enter`; and if they do not work you can use the following keys:

```
Terminal 39: Editing text in Vi
1 h
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

There are some more powerful commands for long-distance jumping.

```
Terminal
```