



Exercises for *Foundations in Data Engineering*, WiSe 18/19

Timo Kersten (kersten@in.tum.de)

<http://db.in.tum.de/teaching/>

Sheet Nr. 01

Exercise 1 Which tools are there in GNU coreutils and for what are they useful? Describe each tool in two sentences.

Coreutils have lots of configuration options, unfortunately often not with coherent command line flags. To learn more about utilities (and other facilities in Linux), use the man page as first source of help. For example type `man join` to learn more about how the join program works.

Exercise 2 Which means does bash offer to combine the coreutils (or any programs for that matter)?

Hint: use `man bash` and read DEFINITIONS, RESERVED WORDS and SHELL GRAMMAR. Also read the beginning of `man 7 pipe` for information on how pipes work in Linux.

A useful bash construct besides redirections (`<`, `>`, `>>` etc.) and pipes (`|`) is process substitution (`<()`). Please read the Wikipedia article if you are not familiar with this.

Exercise 3 For which kinds of data sets are the gnu tools suitable?

Exercise 4 What do we need a pager for even though `cat` can print to the terminal? Name an available pager on Ubuntu.

Exercise 5 Write a regular expression that matches all e-mail addresses in the file `mails_match.txt` and that does not match any of the lines in `mail_dont_match.txt`.

You can easily test your regular expressions using online tools like <https://regex101.com>.

Copy these files to the online tool and try developing your regular expression.

Regular expressions can also be easily used in the command line like this:

```
egrep <your regular expression here> mails_match.txt
```

In a bash in Ubuntu, `egrep` will highlight the matched part per default. This may be helpful for debugging.

To determine whether all the mail addresses matched, you may want to check if the result of your regular expression application has the same number of lines as the original file. That means the following commands must return the same number:

```
wc -l mails_match.txt  
egrep <your regular expression here> mails_match.txt | wc -l
```

Furthermore, the following command should produce no output at all:

```
egrep <your regular expression here> mails_dont_match.txt
```

Exercise 6 It is often useful to divide regular expressions into groups. For example when searching for emails, name and domain are divided by the @ symbol. Groups provide the means within regular expression matching to match name and domain separately.

There are two different types of groups, capturing and non-capturing groups.

Capturing groups are used to store the matched string for later use like replacing. Non-capturing groups help to structure your regular expression.

Use non-capturing groups to find a short expression to match MAC addresses in the file random-mac.txt. Note that grep needs the flag -E to support groups.

Use capturing groups and replace to change the structure of the file names.txt from "First-name Lastname" to "Lastname, Firstname". You may notice that GNU utils are not the right tool for this task. Consider using a scripting language, e.g.:

```
perl -n -e '/.*\/_&&_print_$1'
```

Exercise 7 Convert RGBA color codes in the form "#FF0000FF" (red) to the rgba form "rgb(255, 0, 0, 1.0)" The first three digits indicate how much off the additive light basic colors have to be added to generate the color. The fourth entry determines the alpha value, also known as opacity which ranges from 0 to 1.