

Introduction to **awk** programming

(block course)

Solutions to the exercises



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Michael F. Herbst

michael.herbst@iwr.uni-heidelberg.de

<http://blog.mfhs.eu>

Interdisziplinäres Zentrum für wissenschaftliches Rechnen
Ruprecht-Karls-Universität Heidelberg

15th – 17th August 2016

Contents

Contents	i
Solutions to the exercises	1
Solution to 1.3	1
Solution to 1.4	1
Solution to 2.2	2
Solution to 2.3	2
Solution to 2.4	2
Solution to 2.5	3
Solution to 3.1	3
Solution to 3.3	4
Solution to 3.5	4
Solution to 3.6	5
Solution to 3.7	6
Solution to 3.8	6
Solution to 3.10	7
Solution to 3.11	9
Solution to 4.1	9
Solution to 4.2	10
Solution to 4.5	11
Solution to 5.3	11
Solution to 5.5	12
Solution to 6.3	13
Solution to 6.4	14
Solution to 6.5	14
Solution to 6.7	15
Solution to 6.9	16
Solution to 6.10	16
Solution to 6.11	17
Solution to 6.12	18
Solution to 7.3	19
Solution to 8.3	20
Solution to 9.4	21
Licensing and redistribution	23

Solutions to the exercises

Solution to 1.3

- The commands are

```
$ awk '/hunger/' resources/gutenberg/pg76.txt
```

and

```
$ awk '/hunger/' resources/gutenberg/pg74.txt
```

- One solution for the program `first.awk` is

```
1 # print lines matching hunger
2 /hunger/ {print}
3
4 # print lines matching year
5 /year/ {print}
```

1_first_look/sol/first.awk

which works since if either rule matches the line is printed. Unfortunately it will print the line twice, if both the words “hunger” as well as “year” are contained in the input data.

Solution to 1.4

- The file `resources/testfile` contains

```
1 some words
2 any data
3 some further words
4 somer morer things
5 more other thing
6 even more data
```

resources/testfile

whereas the output of

```
$ awk -f 1_first_look/printprint.awk resources/testfile
```

is

```
1 some_words
2 some_words
3 any_data
4 any_data
5 some_further_words
6 some_further_words
7 somer_morer_things
8 somer_morer_things
```

```

9 more_ other_ thing
10 more_ other_ thing
11 even_ more_ data
12 even_ more_ data

```

in other words, each line is shown twice. This is due to the fact, that the `awk` program `1_first_look/printprint.awk` contains the rule `{print}` twice, which is unconditionally executed. Therefore this printing instruction (which prints the whole record) is executed twice, i.e. the output contains each line twice.

- Now `awk` reads no records (empty file) and hence none of the two rules is executed. Therefore the program produces no output.

Solution to 2.2

The matching part:

- `..` matches any string that contains any two character substring, i.e. any string with two or more letters. This is everything except `g` and the empty string.
- `^..$` matches a string with exactly two characters, i.e. `ab` and `67`.
- `[a-e]` matches any string that contains at least one of the characters `a` to `e`, i.e. `ab` and `7b7`.
- `^.7*$` matches any string which starts with an arbitrary character and then has zero or more `7`s following. This is `g`, `67`, `67777`, `7777` and `77777`.
- `^(.7)*$` matches any string which has zero or more consecutive substrings consisting of an arbitrary character and a `7`. This is `67`, `o7x7g7`, `7777` and the empty string. Note that e.g. `77777` does not match: If we “use” the pattern `.7` three times we get `^.7.7.7$` and `77777` has one character too little to be a match for this.

Solution to 2.3

The crossword:

	<code>a?[3[:space:]]+b?</code>	<code>b[~eaf0-2]</code>
<code>[a-f][0-3]</code>	<code>a3</code>	<code>b3</code>
<code>[:xdigit:]+b+</code>	<code>3b</code>	<code>bb</code>

Solution to 2.4

- `ab*c` or `c$` or just `c`
- `ab+c` or `bc$`
- `^a.*c` or `c$`
- `^_*q` or `q..`
- `^a|w` or `....`

Solution to 2.5

- Regexes for the parts:
 - sign: “[+-]”
 - prefactor: “[01]\.[0-9]*”
 - exponent: “[0-9]+”
- So altogether the scientific numbers need to match:

```
1 ([+-]?)([01]\.[0-9]*)e([+-]?)([0-9]+)
```

where the parenthesis () are only provided to show the individual parts, i.e.

```
1 [+ -]?[01]\.[0-9]*e[+-]?[0-9]+
```

would be valid as well. Executing this on the `digitfile` gives

```
$ awk '/[+-]?[01]\.[0-9]*e[+-]?[0-9]+/ {print}' ↙  
↪resources/digitfile
```

```
1 1.759e+15  
2 1.5e+5da is a scientific number  
3 -1.34e+04
```

- Introducing the fault tolerance implies:
 - We replace the plain requirement for “e” by the bracket expansion “[eEdD]”.
 - Instead of “[01]\.[0-9]*”, we require a number with an optional decimal part, i.e. “[0-9]+(\.[0-9]*)?”

Hence overall

```
1 [+ -]?[0-9]+(\.[0-9]*)?[eEdD][+-]?[0-9]+
```

```
$ awk '/[+-]?[0-9]+(\.[0-9]*)?[eEdD][+-]?[0-9]+/ {print}' ↙  
↪resources/digitfile
```

```
1 1.759e+15  
2 -9.3e-5  
3 19e-5 is not properly formatted either.  
4 1.5e+5da is a scientific number  
5 -1.34e+04
```

Solution to 3.1

This can be achieved using the commandline

```
$ awk '/^free/ { print $2 }' resources/gutenberg/pg1232.txt
```

which prints

```
1 Republic
2 him
3 he
4 sharing
5 distributed
```

Solution to 3.3

The first and second column from the matrix file can be extracted using

```
$ awk '{ print $1 "_" $2 }' resources/matrices/lund_b.mtx
```

which gives

```
1 %%MatrixMarket_matrix
2 %_This
3 147_147
4 1_1
5 2_1
6 8_1
7
8 ...
9
10 146_146
11 147_146
12 147_147
```

If we want to exclude the first two lines (comment lines), we need to run

```
$ awk '/^[^%]' { print $1 "_" $2 }' resources/matrices/lund_b.mtx
```

instead.

Solution to 3.5

One solution to the exercise is

```
1 {
2   res = res "_" $0
3   print res
4 }
```

3_basics/sol/growingconcat.awk

If additionally one wants to get rid of the leading space in each line, one could use the program

```
1 {
2   res = res blank $0
3   blank = "_" # set blank to be a space from here on
4   print res
```

The idea behind this latter script is, that for the first record `blank` and `res` are not defined, i.e. equivalent to the empty string.

Solution to 3.6

First we explain the program:

- The first line of `3_basics/exscript.awk` just causes the current value of the variable `a` to be printed. If this variable is undefined or empty it will print an empty line.
- The second line always sets `num` to the string `"false"` and increases the value of `a`.
- Third line decreases `a` and sets `num` to `"true"` if the record, which is processed contains a digit `0...9`
- In other words if the record contains a digit the value of `a` will overall remain unchanged and `num` is `"true"` before executing line 4.
- Line 4 will just print the value of `num`, so if this line prints `num: false` then the value of `a` is increased.

Now we look at the input.

- The first record is `4`. Here no value resides in `a`, i.e. we print an empty string. Furthermore `num` is set to `"true"` and `a` is updated to `0`. The output of this record is

```
1
2 num: true
```

- Next record is a number as well. We print the `0` from the previous record and the same `num: true`. No change to `a`. The output is

```
1 0
2 num: true
```

- Next record contains no number, so `a` is increased to `1` and `num` is now `"false"`, which yields

```
1 0
2 num: false
```

- Finally we print the increased `a` and increase it further, since `num` is still `"false"`:

```
1 1
2 num: false
```

- and so on

Solution to 3.7

In order to count the number of lines which contain any digit, we can use the script

```
1 /[0-9]/ { c+=1 }
2 END { print c }
```

3_basics/sol/count_numbers.awk

This will provide us with those lines containing *any* kind of number as well, since numbers are obviously made up of digits.

The program

```
1 /[+-]?[01]\.[0-9]*e[+-]?[0-9]+/ { c+=1 }
2 END { print c }
```

3_basics/sol/count_scinumbers.awk

on the other hand counts the number of lines with scientific numbers (in the strict sense).

Solution to 3.8

We compute the column-wise averages using the program

```
1 {
2   count++      # Count of the matrix elements
3   sum1 += $1   # Sum of first column
4   sum2 += $2   # Sum of second column
5   sum3 += $3   # Sum of third column
6 }
7
8 END {
9   # Compute averages and print:
10  print "Average_col1:_" sum1/count
11  print "Average_col2:_" sum2/count
12  print "Average_col3:_" sum3/count
13 }
```

3_basics/sol/mtx_averages.awk

This results in

```
$ awk -f 3_basics/sol/mtx_averages.awk resources/matrices/3.mtx
```

```
1 Average_col1:_1.90909
2 Average_col2:_1.90909
3 Average_col3:_2.45455
```

and

```
$ awk -f 3_basics/sol/mtx_averages.awk ↙
↪resources/matrices/lund_b.mtx
```



```
1 Average_col1: 79.1026
2 Average_col2: 68.2924
3 Average_col3: 150.228
```

and

```
$ awk -f 3_basics/sol/mtx_averages.awk ↙
↪resources/matrices/bcsstm01.mtx
```

```
1 Average_col1: 24.48
2 Average_col2: 24.48
3 Average_col3: 64.96
```

If one wants to make sure to skip the first few comment lines, one can use the program

```
1 /^[^%]/ {
2     count++           # Count of the matrix elements
3     sum1 += $1        # Sum of first column
4     sum2 += $2        # Sum of second column
5     sum3 += $3        # Sum of third column
6 }
7
8 END {
9     # Compute averages and print:
10    print "Average_col1:" sum1/count
11    print "Average_col2:" sum2/count
12    print "Average_col3:" sum3/count
13 }
```

3_basics/sol/mtx_averages_skip.awk

instead, where a guarding regular expression pattern makes sure that only non-comment lines are included in the average.

Solution to 3.10

- The program

```
1 # check if we have a comment. If not increase the line number
2 # and flag as a nocomment record
3 /^[^%]/ {
4     linenr +=1
5     nocomment=1
6 }
7
8 # Extract the number of entries and store them
9 linenr == 1 && nocomment {
10    nentries=$3
11 }
12
13 # Increase the count of actual entries,
14 # since this an explicitly provided entry
```

```

15 linenr > 1 && nocoment {
16     actualentries++
17 }
18 END {
19     print "Expected_entries:#####" nentries
20     print "Actual_entries:#####" actualentries
21 }

```

3_basics/sol/mtx_check_entry_count.awk

prints both the counted and the expected number of non-zero entries in the `mtx` file.

- Both values for the sparsity ratio are printed by

```

1 # check if we have a comment. If not increase the line number
2 # and flag as a nocoment record
3 /^[^%]/ {
4     linenr +=1
5     nocoment=1
6 }
7
8 # Extract the number of nonzeros and store them
9 # Compute the number of rows times columns
10 linenr == 1 && nocoment {
11     rows = $1
12     cols = $2
13     nentries=$3
14     total = rows*cols
15 }
16
17 # Increase the count of actual entries,
18 # since this an explicitly provided entry
19 linenr > 1 && nocoment {
20     actualentries++
21 }
22
23
24 END {
25     print "Sparsity_ratio:#####" (total-nentries)/total
26     print "Actual_sparsity_ratio:#" (total-actualentries)/total
27 }

```

3_basics/sol/mtx_sparsity_ratio.awk

- The elementwise square is computed by the program

```

1 # Copy all comments:
2 /^%/ { print $0 }
3
4 # If no comment and we have passed the first line:
5 /^[^%]/ && passedfirst == 1 {
6     print $1 "□" $2 "□" ($3*$3)
7 }
8
9 # Copy the first non-comment line verbatim:

```

```

10 /^[^%]/ && passedfirst != 1 {
11     passedfirst=1
12     print
13 }

```

3_basics/sol/mtx_square_elements.awk

Solution to 3.11

One possible way to extract the excited states is:

```

1 #!/usr/bin/awk -f
2
3 # We use the state variable inside_block to keep track whether
4 # we are inside or outside an excited states block
5 # It's default value is 0, i.e. outside
6
7 # whenever we encounter the " Excited state ", we
8 # change the flag to indicate that we are inside the table.
9 # also we store the state number, which sits in the third field
10 /^ *Excited state / { inside_block=1; state_number=$3 }
11
12 # if we find the "Term symbol" line inside the block, we store
13 # the term symbol which sits in $3 $4 and $5
14 inside_block==1 && /^ *Term symbol/ { term_symbol=$3 "␣" $4 "␣" ↙
    ↪$5 }
15
16 # if we find the "Excitation energy" line, we store the ↙
    ↪excitation energy
17 # and print the table, since we do not care about the rest of the
18 # block. Next we reset the inside_block flag for the next block ↙
    ↪to come.
19 inside_block==1 && /^ *Excitation energy/ {
20     excitation_energy=$3
21
22     # print the data tab-separated (for analysis with e.g. cut)
23     print state_number "\t" term_symbol "\t" excitation_energy
24
25     inside_block=0
26 }

```

3_basics/ex_extract_states.awk

Solution to 4.1

The following program implements one way to print duplicated words in a text:

```

1 #!/usr/bin/awk -f
2 # change the record separator to anything which is not
3 # an alphanumeric (we consider a different word to start
4 # at each alphanumeric character)
5 BEGIN { RS="[^[:alnum:]]+" }

```

```

6 # now each word is a separate record
7
8 $0 == prev { print prev }
9 { prev = $0 }

```

4_parsing_input/sol/duplicated_words.awk

Solution to 4.2

- The final balance is printed by

```

1 #!/usr/bin/awk -f
2
3 # Change field separator:
4 BEGIN { FS = "," }
5
6 # Extract starting balance
7 /^[#] Starting balance/ { balance = $2 }
8
9 # Once in the transfer block, adjust balance:
10 /^[^#]/ { balance += $2 }
11
12 # Print the final balance:
13 END { print "Final balance:" balance }

```

4_parsing_input/sol/csv_balance.awk

- A balance column is appended by

```

1 #!/usr/bin/awk -f
2
3 # Change field separator:
4 BEGIN { FS = "," }
5
6 # Extract starting balance
7 /^[#] Starting balance/ { balance = $2 }
8
9 # Print all comment lines verbatim:
10 /^[#]/
11
12 # Once in the transfer block, adjust balance and append a ↵
13   ↵column
14 /^[^#]/ {
15     balance += $2
16     print $0 "," balance
17 }

```

4_parsing_input/sol/csv_balance_append_column.awk

Solution to 4.5

One solution program to add all scientific numbers which occur in some input is:

```
1 #!/usr/bin/awk -f
2
3 BEGIN {
4   # Pattern for the Sign:
5   sign="[+-]?"
6
7   # Pattern for an integer
8   intp="[0-9]+"
9
10  # For a float we may additionally have
11  fpe="(\.[0-9]+)?"
12
13  # The optional exponent
14  expe="(e[+-][0-9]+)?"
15
16  # Build the pattern:
17  FPAT=sign intp fpe expe
18 }
19
20 # Assume that we have no more than 5 numbers in each line
21 # (which is true for the digitsfile)
22 { c += ($1 + $2 + $3 + $4 + $5) }
23 END { print "The sum is " c }
```

4_parsing_input/sol/add.digits.awk

Solution to 5.3

- One way to achieve the unfolding is to do default input parsing, but to print each field on a different line, e.g.

```
$ awk 'BEGIN { OFS="\n" }; { $1=$1; print }' ✓
↪resources/testfile
```

The other option is to treat each word as a separate record, i.e.

```
$ awk 'BEGIN { RS="[:space:]+" }; { print }' ✓
↪resources/testfile
```

- Changing the separator character in a csv file from comma to semicolon can be achieved by the simple commandline

```
$ awk 'BEGIN {OFS=";"; FS=","}; {$1=$1; print }' ✓
↪resources/data/money.csv
```

which sets FS and OFS appropriately and then triggers a rebuild of the `$0` variable.

Solution to 5.5

One solution to print the average measurement value and to exclude the erroneous apparatus 3 explicitly is

```
1 #!/usr/bin/awk -f
2
3 BEGIN {
4     # make both space and : field separators
5     FS="[:]+"
6
7     # Alternatively we can use FPAT to describe
8     # the numbers that we expect:
9     #
10    # FPAT="[0-9]+|-[0-9]\\.[0-9]+"
11    #
12    # Problem is that this does not work for numbers
13    # in the scientific format like that.
14    # One would need to add another alternation.
15
16    print("##|average")
17    print("----+-----")
18 }
19
20 # only process if the line is no comment line
21 $0 !~ /^#/ {
22     # $1 is the apparatus count
23     # $2 to $8 is the values
24
25     # compute average and add to total sum
26     sum = ($2+$3+$4+$5+$6+$7+$8)
27     avg = sum/7.
28
29     # print avg
30     printf("%2d|%.4f\n",$1,avg)
31 }
32
33 # Apparatus 3 is a little off, so exclude it explicitly
34 $1 != 3 {
35     totsum+= sum
36     totcount+=7
37 }
38
39 # Print a note about this:
40 $1 == 3 { print "|||Note:Not included in total sum" }
41
42 # print results:
43 END { printf("\ntotal avg: %.4f\n",totsum/totcount) }
```

5_printing_output/sol/analysis.awk

Solution to 6.3

If one wants to use a range pattern, this can be done using the program

```
1 #!/usr/bin/awk -f
2
3 # The chapter to extract, here the first
4 BEGIN { v=1 }
5
6 # The range: From this until the next.
7 $1 == "CHAPTER" && $2 == v, $1 == "CHAPTER" && $2 == (v+1)
    6_patterns_actions_variables/sol/extract_chapter.awk
```

Running this like

```
$ 6_patterns_actions_variables/sol/extract_chapter.awk ✓
  ↪resources/gutenberg/pg161.txt
```

gives

```
1 CHAPTER_1
2
3
4 The_ family_ of_ Dashwood_ had_ long_ been_ settled_ in_ Sussex.  Their_ ✓
  ↪estate
5 was_ large,  and_ their_ residence_ was_ at_ Norland_ Park,  in_ the_ ✓
  ↪centre_ of
6
7 ...
8
9 having_ much_ of_ her_ sense,  she_ did_ not,  at_ thirteen,  bid_ fair_ to_ ✓
  ↪equal
10 her_ sisters_ at_ a_ more_ advanced_ period_ of_ life.
11
12
13
14 CHAPTER_2
```

In order to avoid the chapter heading of the next chapter to be printed, one could store the chapter number instead:

```
1 #!/usr/bin/awk -f
2 BEGIN { v=1 }
3
4 # remember chapter number
5 /^CHAPTER [0-9]+/ { chapter = $2 }
6
7 # Print the chapter if it is the right one
8 chapter == v
    6_patterns_actions_variables/sol/extract_chapter_state.awk
```

Solution to 6.4

One solution is to count the number of lines inside the Davidson range:

```
1 #!/usr/bin/awk -f
2
3 # We know that the iteration count increases by one
4 # for each extra line we find in the Davidson block
5 # There are 7 lines containing no iterations
6 # (ie the headings, the guess and the summary)
7 # So we count all lines between "Starting Davidson"
8 # and "Davidson Summary" and subtract 7 to get the
9 # number of iterations.
10 /^ Starting Davidson \.\.\./, /^ Davidson Summary:/ { count+=1 }
11 /^ Davidson Summary:/ {
12     # print count and reset
13     print count-7
14     count=0
15 }
```

6_patterns_actions_variables/sol/extract_davidson.awk

Solution to 6.5

If one wants to automatically exclude the instrument based on an upper threshold, one could use the program

```
1 #!/usr/bin/awk -f
2
3 BEGIN {
4     # The upper threshold to include a value
5     thresh_upper = -0.05
6
7     # make both space and : field separators
8     FS="[_:]+"
9
10    print("##_|_average")
11    print("----+-----")
12 }
13
14 # Skip comment lines:
15 /^#/ { next }
16
17 # $1 is the apparatus count
18 # $2 to $8 is the values
19
20 {
21     # Compute the average:
22     sum = ($2+$3+$4+$5+$6+$7+$8)
23     avg = sum/7.
24
25     # if the average is larger than upper threshold,
26     # the apparatus is off and we skip the rest
```



```

27 if (avg >= -0.05) {
28     printf("%2d|%.4f>%.4f=>Not\n", $1, avg, thresh_upper)
29     next
30 }
31
32 # All the records that made it here
33 # should be included:
34 printf("%2d|%.4f\n", $1, avg)
35
36 totsum+= sum
37 totcount+=7
38 }
39
40 END { printf("\ntotal avg: %.4f\n", totsum/totcount) }

```

6_patterns_actions_variables/sol/analysis_automatic.awk

Solution to 6.7

- A possible factorial program is

```

1 #!/usr/bin/awk -f
2 {
3     n+=$1
4     res=1
5     while(n>1) {
6         res=res*n
7         --n
8     }
9     print res
10 }

```

6_patterns_actions_variables/sol/factorial.awk

- A couple of examples:

```

$ echo -e "20\n50\n100" |
  ↪6_patterns_actions_variables/sol/factorial.awk

```

gives

```

1 2432902008176640000
2 3041409320171337557636696640...832057064836514787179557289984
3 9332621544394410218832560610...311236641477561877016501813248

```

So `awk` is able to do integer arithmetic up to the point that it allows to calculate $100!$ using an extremely naive algorithm!

Solution to 6.9

Just replace the `while` by a `for` loop:

```
1 #!/usr/bin/awk -f
2 {
3     res=1
4     for(n=+$1; n>0; --n) {
5         res=res*n
6     }
7     print res
8 }
```

6_patterns_actions_variables/sol/factorial_for.awk

Solution to 6.10

We generalise the program by using a `for`-loop over fields and make the code cleaner using an `if`-statement.

```
1 #!/usr/bin/awk -f
2
3 BEGIN {
4     # The upper threshold to include a value
5     thresh_upper = -0.05
6
7     # make both space and : field separators
8     FS="[:]+"
9
10    print("##|average")
11    print("---+-----")
12 }
13
14 # Skip comment lines:
15 /^#/ { next }
16
17 # $1 is the apparatus count
18 # from $2 onwards are the values
19
20 # Compute the average:
21 {
22     # Accumulate the sum:
23     sum=0
24     for (i=2;i<=NF;++i) {
25         sum+=$i
26     }
27     avg = sum/(NF-1)
28
29     # check if the average is larger than
30     # upper threshold, if yes then apparatus
31     # is off and we skip the rest, else we
32     # print and add to the total
33     if (avg > thresh_upper) {
```

```

34     printf("%2d|%.4f>%.4f=>Not\n", $1, avg, thresh_upper)
35     ↪included\n", $1, avg, thresh_upper)
36     next
37 } else {
38     printf("%2d|%.4f\n", $1, avg)
39     totsum+= sum
40     totcount+=(NF-1)
41 }
42 }
43 END { printf("\ntotal avg: %.4f\n", totsum/totcount) }

```

6_patterns_actions_variables/sol/analysis_general.awk

Solution to 6.11

The following script checks whether the first field of each record is a prime number.

```

1 #!/usr/bin/awk -f
2 {
3
4     isprime=1
5     n=+$1
6     for (i=2; i*i < n; ++i) {
7         if (n % i == 0) {
8             isprime=0
9             break
10        }
11    }
12
13    if (isprime) {
14        printf("%d is prime\n", n)
15    } else {
16        printf("Smallest divisor of %d is %d\n", n, i)
17    }
18 }

```

6_patterns_actions_variables/sol/is_prime.awk

The output for

```
$ echo -e "101\n1001" | 6_patterns_actions_variables/ex_break.awk
```

is

```

1 101 is prime
2 Smallest divisor of 1001 is 7

```

Solution to 6.12

One solution to find the number of times “a” and “e” occur in a book is

```
1 #!/usr/bin/awk -f
2
3 BEGIN { FS="" }
4
5 /\*\*\* START OF THIS PROJECT GUTENBERG EBOOK [ A-Z,-.]+ \*\*\*/ {
6     # Flag that we are inside, but do not do statics on this record
7     inside_book=1
8     next
9 }
10
11 /\*\*\* END OF THIS PROJECT GUTENBERG EBOOK [ A-Z,-.]+ \*\*\*/ {
12     # We are at the end of the book => quit awk program
13     exit
14 }
15
16 # if we are inside:
17 inside_book {
18     # Ignore case such that both upper
19     # and lower case characters are counted
20     IGNORECASE=1
21
22     for (i=1; i<=NF; ++i) {
23         # Use a regex here, since == operator
24         # is not affected by IGNORECASE
25         if ($i ~ /a/) {
26             account++
27         } else if ($i ~ /e/) {
28             ecount++
29         }
30         charcount++
31     }
32
33     # Unset IGNORECASE, since the regex patters above are
34     # case sensitive.
35     IGNORECASE=0
36 }
37
38 # Print final results:
39 END {
40     printf("total_␣%8d\n", charcount)
41     printf("a_␣␣␣␣␣␣%8d_␣␣(%6.2f%%)\n", account, account/charcount*100)
42     printf("e_␣␣␣␣␣␣%8d_␣␣(%6.2f%%)\n", ecount, ecount/charcount*100)
43 }
```

6_patterns_actions_variables/sol/gutenberg_character_statistics.awk

Solution to 7.3

One solution, which also excludes all whitespace characters when performing the character counting, is

```
1 #!/usr/bin/awk -f
2
3 BEGIN { FS="" }
4
5 /\*\*\* START OF THIS PROJECT GUTENBERG EBOOK [ A-Z,-.]+ \*\*\*/ {
6   # Flag that we are inside, but do not do statics on this record
7   inside_book=1
8   next
9 }
10
11 /\*\*\* END OF THIS PROJECT GUTENBERG EBOOK [ A-Z,-.]+ \*\*\*/ {
12   # We are at the end of the book => quit awk program
13   exit
14 }
15
16 inside_book {
17   for (i=1; i<=NF; ++i) {
18     # Ignore those characters which are space characters:
19     # Note: Not strictly speaking required for the exercise,
20     # but gives a nicer result in the end.
21     if ($i ~ /[[:space:]]/) continue
22
23     # Increase count for character and total count
24     #
25     # one could also use
26     #   count[tolower($i)]
27     # in order to map each character to its lower-case
28     # equivalent and make a count over this instead.
29     count[$i]++
30     charcount++
31   }
32 }
33
34 # Print final results:
35 END {
36   printf("total_␣%8d_✓\n", charcount, charcount/charcount*100)
37   print("-----")
38
39   for (c in count) {
40     printf("%-5s_␣%8d_✓\n", c, count[c], count[c]/charcount*100)
41     ↪ (%6.2f%%)\n", c, count[c], count[c]/charcount*100)
42   }
43 }
```

7.arrays/sol/gutenberg_character_statistics.awk

Solution to 8.3

If we allow ourselves to use the usual control structures one could find the maximum and absolute maximum like this

```
1 #!/usr/bin/awk -f
2
3 # The usual abs function
4 function abs(a) {
5     if (a<0) return -a
6     return +a
7 }
8
9 # Initialise max and absmax:
10 NR == 1 {
11     max = $1
12     absmax = abs($1)
13 }
14
15 # Loop over each field (number) and update
16 # max and absmax if necessary
17 {
18     for(i=1;i<=NF;++i) {
19         if ($i > max) {
20             max = $i
21         }
22         if (abs($i) > absmax) {
23             absmax=abs($i)
24         }
25     }
26 }
27
28 END {
29     print "max:UUUU" max
30     print "absmax:U" absmax
31 }
```

8.functions/sol/max_element.long.awk

Alternatively, we can change the range separator and use `awk`'s implicit loop over records to achieve the same thing in less lines of code and without a single control structure:

```
1 #!/usr/bin/awk -f
2
3 # The usual abs function
4 function abs(a) {
5     if (a<0) return -a
6     return +a
7 }
8
9 # Change record separator to repeated space chars
10 # so each field of the matrix becomes a record on its own.
11 BEGIN { RS="[[[:space:]]+" }
12
13 # Initialise max and absmax with first record:
```

```

14 NR == 1 {
15     max = +$0
16     absmax = abs($0)
17     next
18 }
19
20 # For all other record, determine if max or absmax:
21 +$0 > max { max = +$0 }
22 abs($0) > absmax { absmax = abs($0) }
23
24 END {
25     print "max:␣␣␣␣" max
26     print "absmax:␣" absmax
27 }

```

8_functions/sol/max_element.awk

Solution to 9.4

- `wc -w` is equivalent to

```

1 #!/usr/bin/awk -f
2 # Split into a new record at multiple occurrences of space
3 # characters. Then just print the record count.
4 BEGIN { RS="[:space:]+" }
5 END { print NR }

```

9_practical_programs/sol/wc_w.awk

- `uniq -c` we can implement like

```

1 #!/usr/bin/awk -f
2
3 # Initialise buffer to be the first record:
4 NR == 1 { buffer=$0 }
5
6 # If repeated occurrence increase count:
7 buffer == $0 { count++ }
8
9 # Else print the record we had in the buffer
10 # and reset counter and buffer
11 buffer != $0 {
12     printf("%5d␣%s\n",count,buffer)
13     buffer=$0
14     count=1
15 }
16
17 # Print what is left in the buffer
18 END {
19     printf("%5d␣%s\n",count,buffer)
20 }

```

9_practical_programs/sol/uniq_c.awk

- `sort` is implemented using `awk`'s `asort`:

```
1 #!/usr/bin/awk -f
2
3 # Append all input lines to a buffer array
4 { buffer[NR] = $0 }
5
6 # In the end sort using asort and print in order
7 END {
8     nr = asort(buffer)
9     for (i=1; i<=nr; ++i) {
10         print(buffer[i])
11     }
12 }
```

9-practical-programs/sol/sort.awk

- `egrep` can be mimicked using a surrounding shell script with inline `awk` code:

```
1 #!/bin/sh
2 # Store the regex (first argument to script)
3 regex=$1
4 shift
5
6 # Call awk and use DOUBLE quotes to insert the regex
7 # inside an awk pattern and pass the remaining
8 # arguments to the scripts to awk itself (as files)
9 #
10 # Whenever that regex pattern matches the default print
11 # action is executed (exactly like egrep does it)
12 awk "/$regex/" $@
```

9-practical-programs/sol/egrep.sh

For more details, how the shell command `shift` works and what the shell variables `$1` and `$@` mean, see chapter 3.2.1 and 4.6 of the lecture notes to the “advanced bash scripting” course¹.

¹Available from <http://blog.mfhs.eu/teaching/advanced-bash-scripting-2015/>.

Licensing and redistribution

This work is licensed under the Creative Commons Attribution-ShareAlike 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/4.0/>.



An electronic version of this document is available from <http://blog.mfhs.eu/teaching/introduction-to-awk-programming-2016/>. If you use any part of my work, please include a reference to this URL along with my name and email address.