# Assignment 6
## Exercise I

**Download the book „The Adventures of Tom Sawyer" from http://www.gutenberg.org/ebooks/74 (utf-8 format) and store it locally.**

**SOLVE**

Firstly, I checked that there is no files in ls. Then I opened a new txt file with touch. I copied the given text into this file.



And I saved the given text to the file.

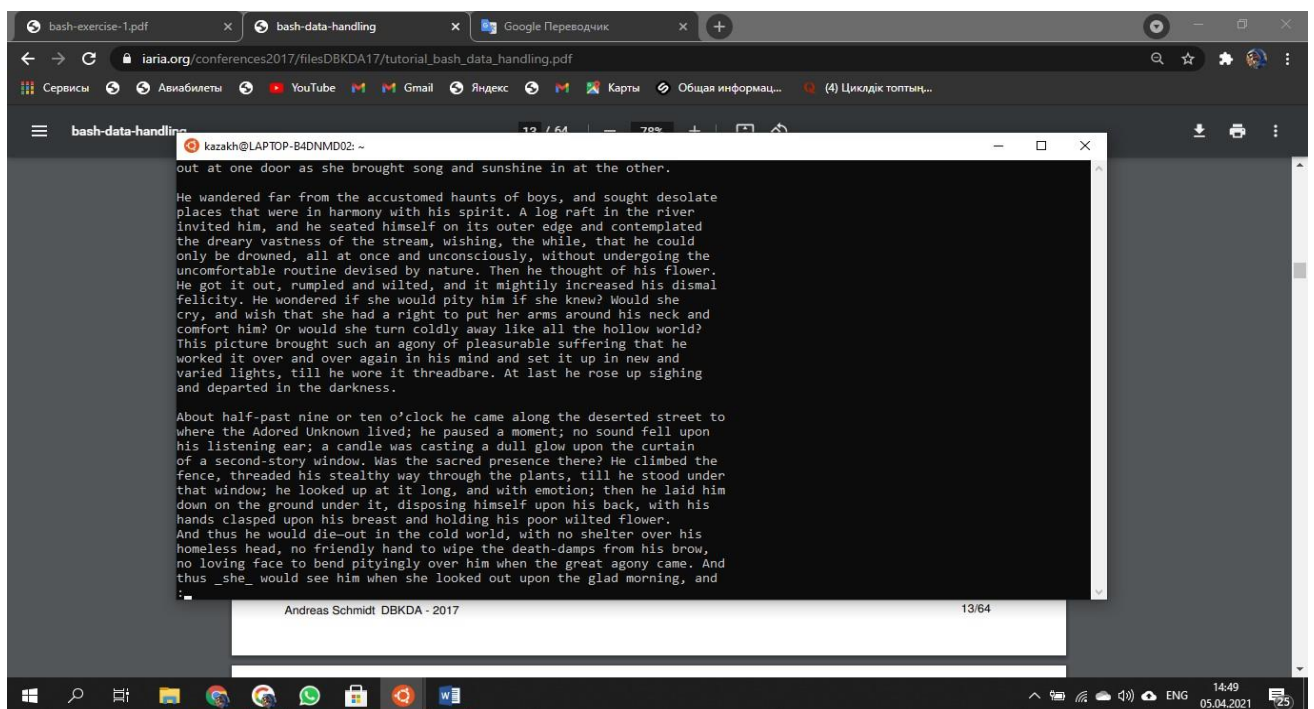**2. For cygwin users only: Convert file to Unix Format with command: dos2unix.exe**

**SOLVE:** I'm done


**3. Browse through the first pages of the book using less (for help type man less)**

**a) Go to line 1234 of the file. What ist the third word?**

-------------------------------------------------------------------------------------------------------------------------

**SOLVE**

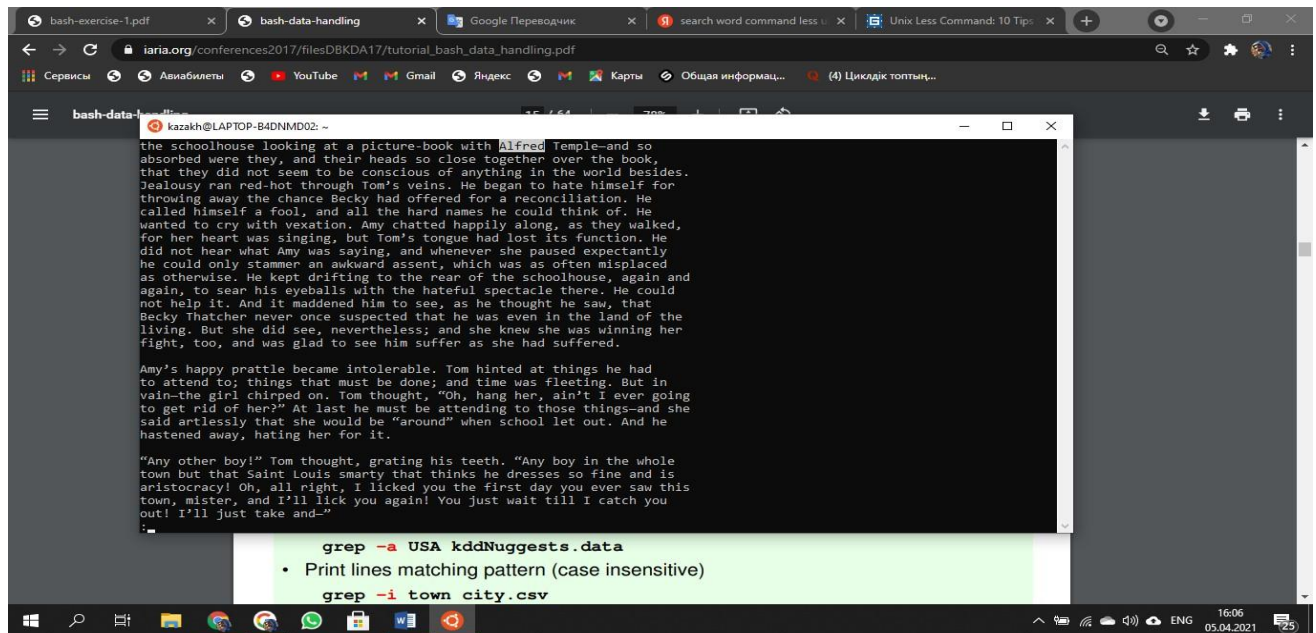Firstly, I entered command **less ass6.txt,** then entered **1234g** and went to a given line.

I can read the third word without any code.) Third word is one.

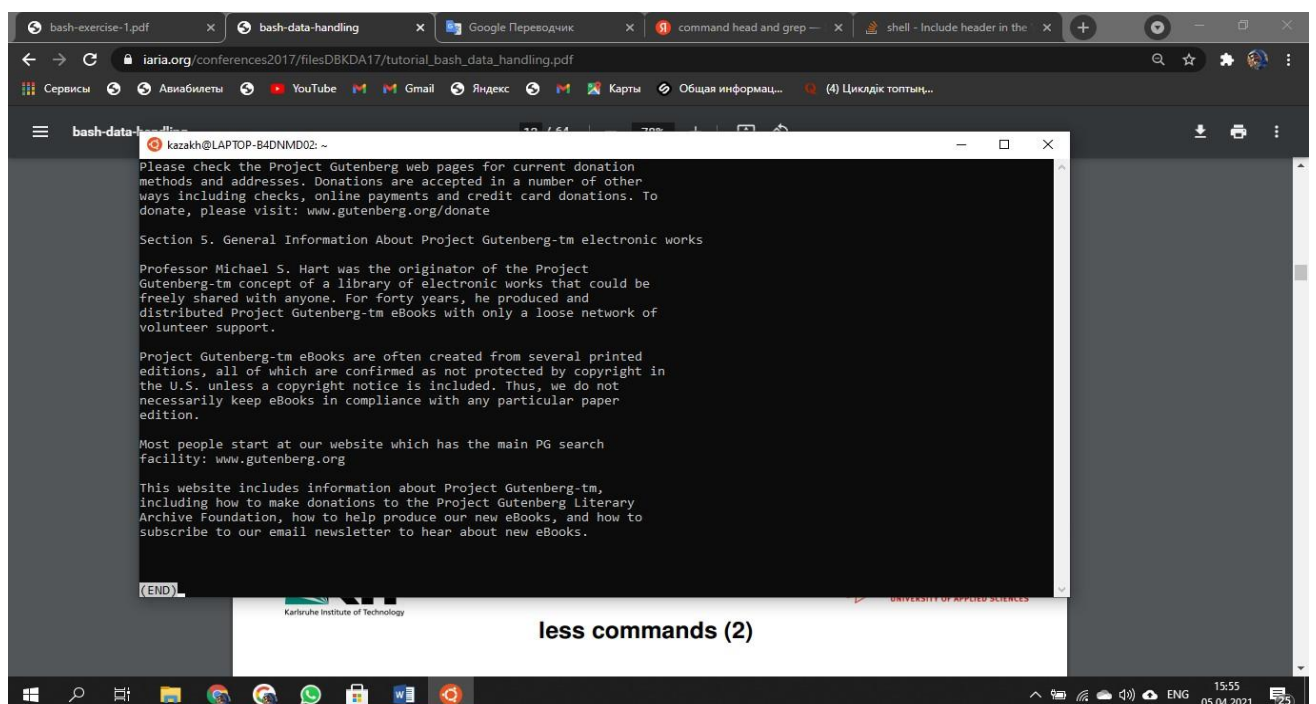**b) Search for the first apperance of Alfred (what is his last name?)**

_____

**SOLVE**

I used the command "less" and entered "/Alfred". And found first line and his last name. His last name is Temple.



**c) Go to the end of the document.**

**SOLVE**

I used the "less" command again and entered ">". This command takes us to the end of the text.

## d) Scroll backward 5 pages.

**SOLVE**

Going to the end of the text, I gave the "b" command 5 times.



## e) Show only lines containing the word „Tom".

**SOLVE**

I used "grep" for search all lines. I entered "grep Tom ass6.txt".

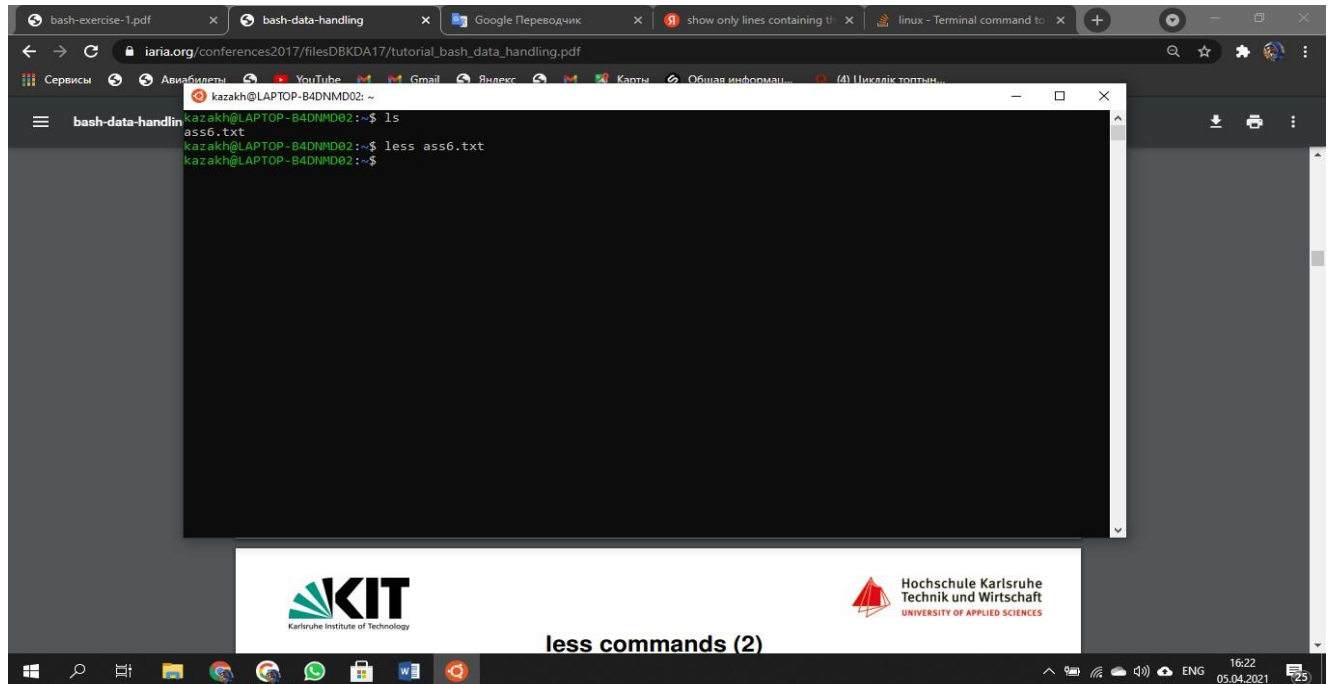**f) Quit less**

**SOLVE**

It's easy. I entered only "q" and I got out of there.



**4. How many chapters has the book (hint use grep with –a option)?**

**SOLVE**

I use "grep" command and entered "grep –ia ass6.txt". ~i~ is ignore case and a – matching lines a binary file. Here are all the chapters and their number is 35

**5. Count the number of empty lines in the file.** ----------------------------------------------------
-------------------------------------------------------------------

**SOLVE**

I use grep and entered "grep -E -wc '^$' ass6.txt". Here -E receives only the marked text. And -wc gives their number. So, empty lines are 2320.



**6. Execute grep Tom and grep –o Tom . What is the difference?**

-----------------------------------------------------------------------------------------------------------------

**SOLVE**

I search for sentences with the word "Tom" in the "grep Tom file" command.

In the "grep -o Tom file" command, words other than "Tom" were cut out.

## 7. How often does the names „Tom" and „Huck" appears in the book?

--------------------------------------------------------------------------------------------------------------------

**SOLVE**

I use grep and entered "grep –wc name file". –wc counts the given word.

So, Tom-785 and Huck-251.

```
kazakh@LAPTOP-B4DNMD02:~$ grep -wc "Tom" ass6.txt
785
kazakh@LAPTOP-B4DNMD02:~$ grep -wc "Huck" ass6.txt
251
kazakh@LAPTOP-B4DNMD02:~$
```

# Assignment 6

## Exercise 2

**1. Download the book 'The Adventures of Tom Sawyer' from http://www.gutenberg.org/files/74/74-0.txt and store it locally in a file called 'The-Adventures-of-Tom-Sawyer.txt.'**

**SOLVE**

I opened a new txt file with touch. I copied the given text into this file. And I named it "The-Adventures-of-Tom-Sawyer.txt".

**2. Count the words and lines in the book 'The-Adventures-of-Tom-Sawyer.txt'.**

**SOLVE**

**I entered "wc –l file.txt" for lines. There are 9255 lines.**

**I entered "wc –w file.txt" for words. There are 73872 words.**

## 3. What does the following command perform?

egrep -o '[A-Za-z]+' The-Adventures-of-Tom-Sawyer.txt

## SOLVE

This code divides words in the text with the symbols [A-Z, a-z] into each line. And -o removes the rest.

**4. Translate all words of 'The-Adventures-of-Tom-Sawyer.txt' into lowercase using the command tr.**

**SOLVE**

I used tr and I wrote all the words in lower case using this code.

tr 'A-Z' 'a-z' < The-Adventures-of-Tom-Sawyer.txt

**5. Count, how often each word in this book appears (hint: use sort, uniq).**

**SOLVE**

I tried to do it based on the solution of the previous task.

I sorted the words, then counted them.



```
kazakh@LAPTOP-B4DNMD02:~$ egrep -o '[A-Za-z]+' The-Adventures-of-Tom-Sawyer.txt | sort | uniq -c
    94 A
     2 ABOUT
     1 ACTUAL
     3 ADVENTURES
     1 AFTER
     2 AGREE
     1 AGREEMENT
     1 ALABAMA
     1 AND
     3 ANY
     1 ANYTHING
     2 AS
     2 ASCII
     2 AT
     1 AUTHOR
     8 About
     1 Academy
     1 Acrobatics
     2 Act
     1 Additional
     1 Adored
     1 Advantages
     3 Adventures
     1 Afeard
     1 Afraid
    16 After
     1 Against
     1 Agreed
     5 Ah
```
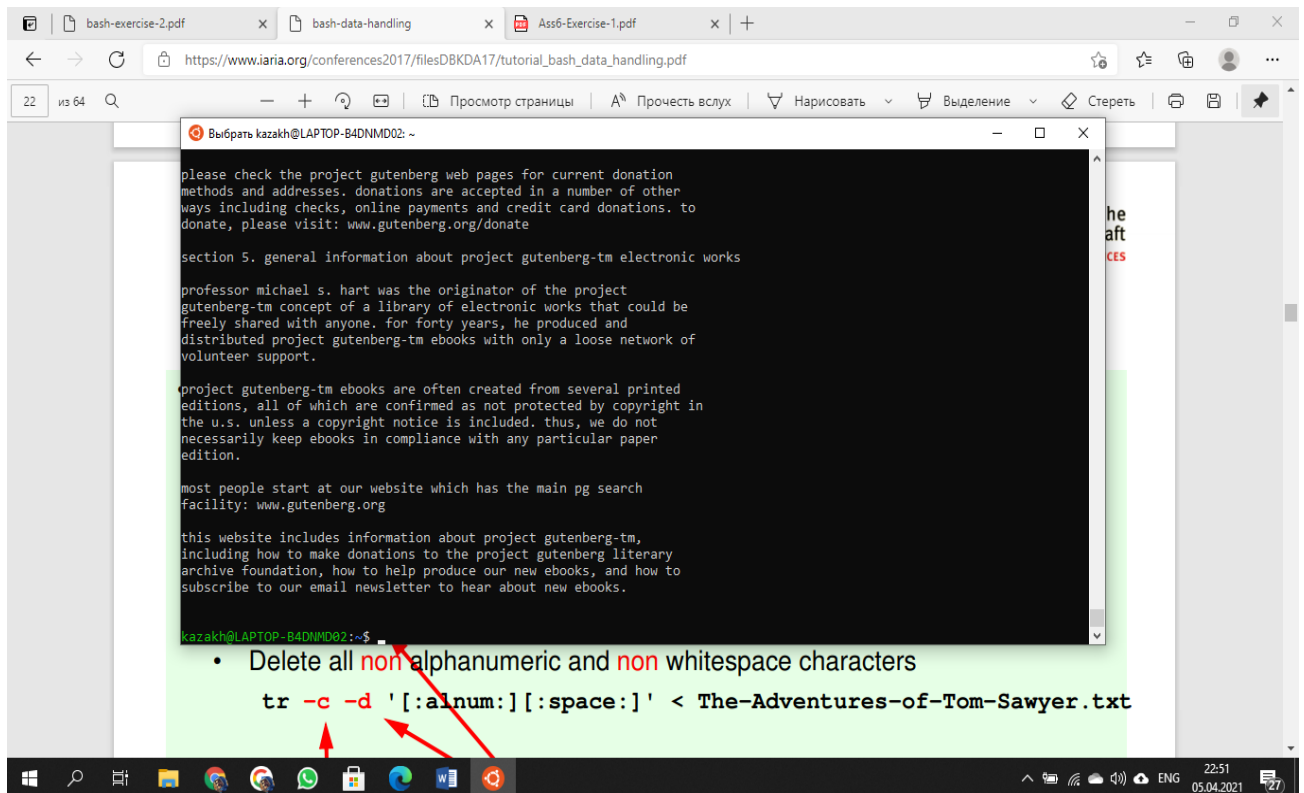
**6. Order the result, starting with the word with the highest frequency. Which word is it?**

**SOLVE**

I added sort -nr to my code. This word is "the".

## 7. Write all the above steps in one statement (using pipes)

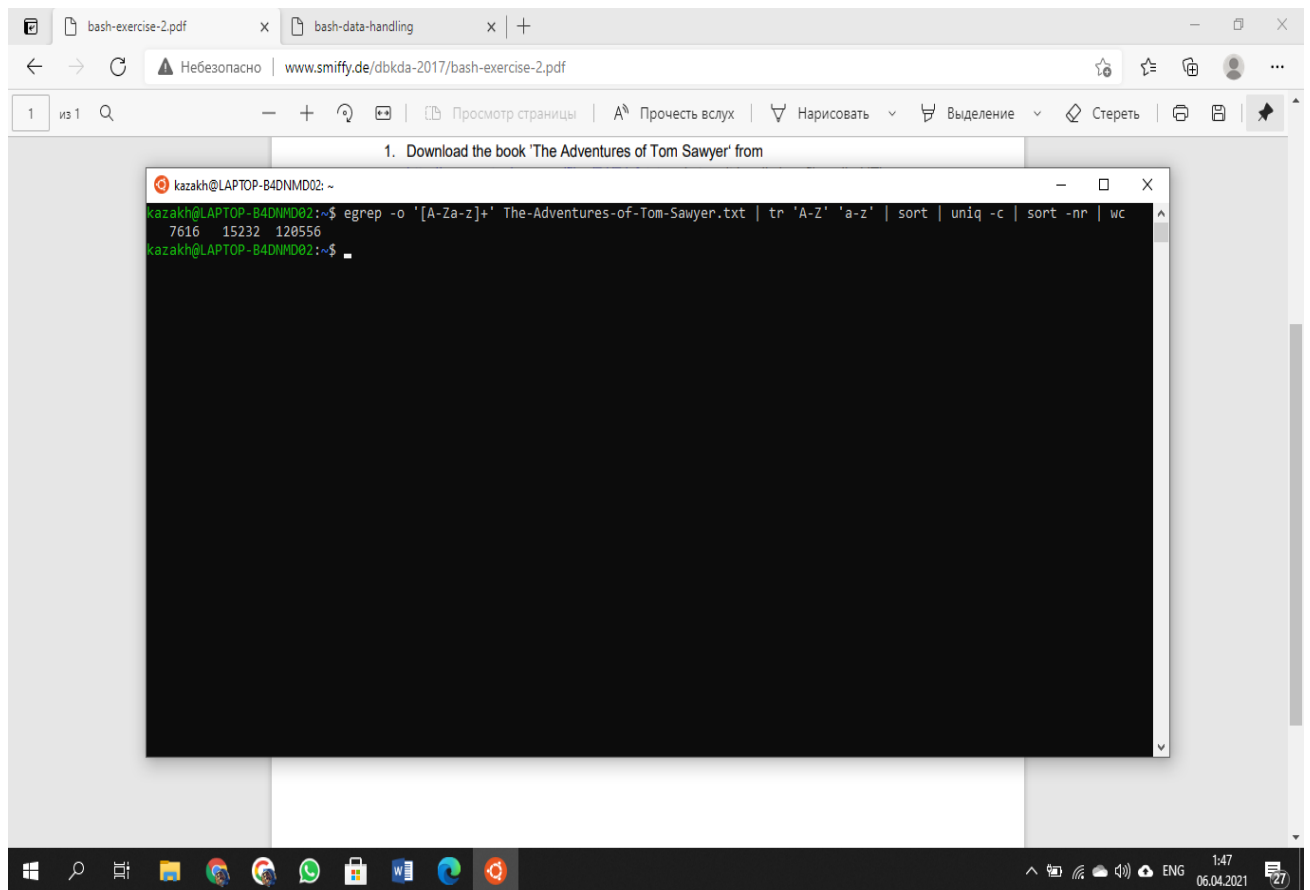**SOLVE**

I wrote one code using all the commands. Lines – 7616 and words – 15232.

I used pipes | |.

**8. Compare the result with the result from the following book: http://www.gutenberg.org/files/2701/2701-0.txt. At which positions (rank) appear the first book specific words?**

**SOLVE**

I use "comm". **comm** is compare two sorted files line by line.
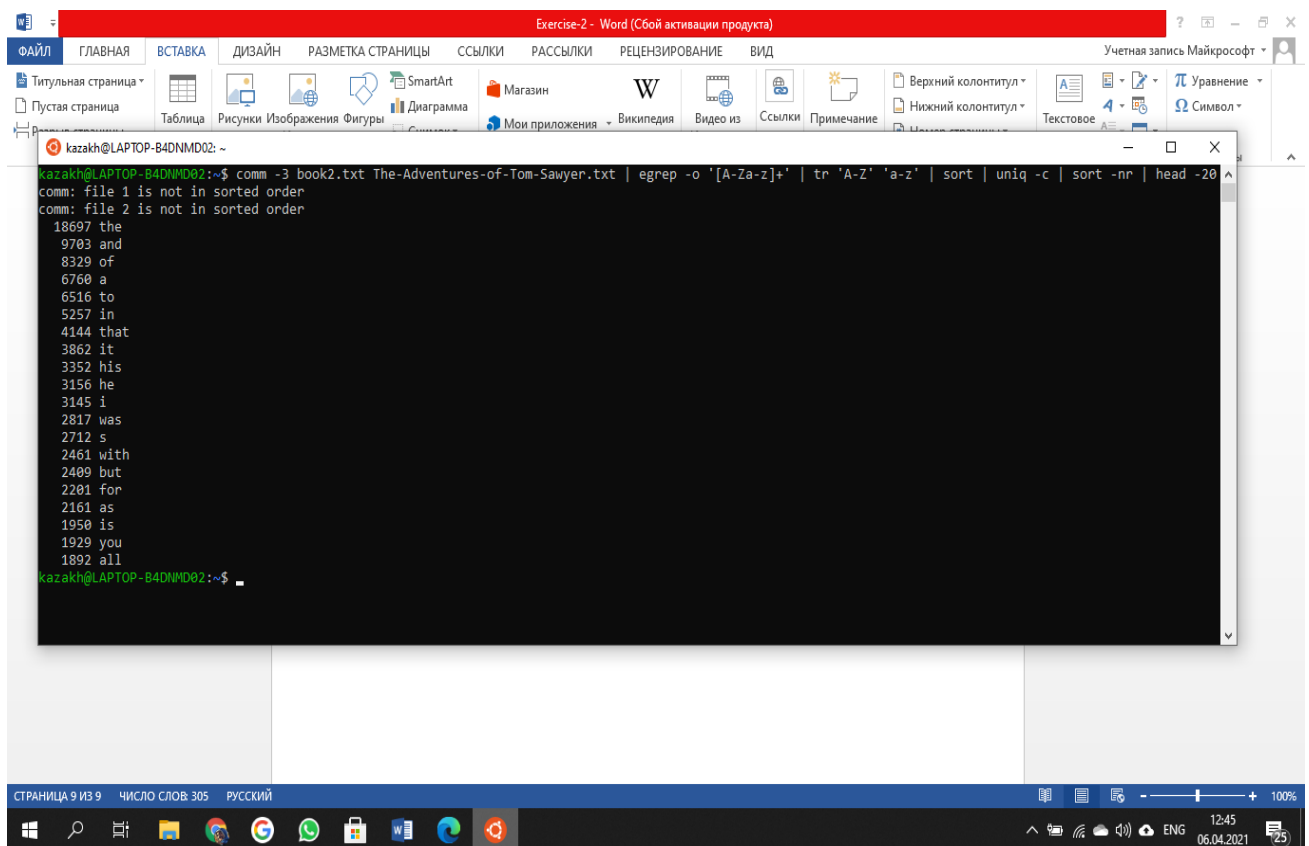
My code: **comm file1 file2**

First book specific words in first column.

**9. Compare the 20 most frequent words of each book. How many are in common? (hint: use head, cut, comm).**

**SOLVE**

I chose the lines in both books with comm. I put each word on a line with egrep. I ignored the indexes with tr. Sorted, counted, sorted by frequency and print the 20 most frequent words.
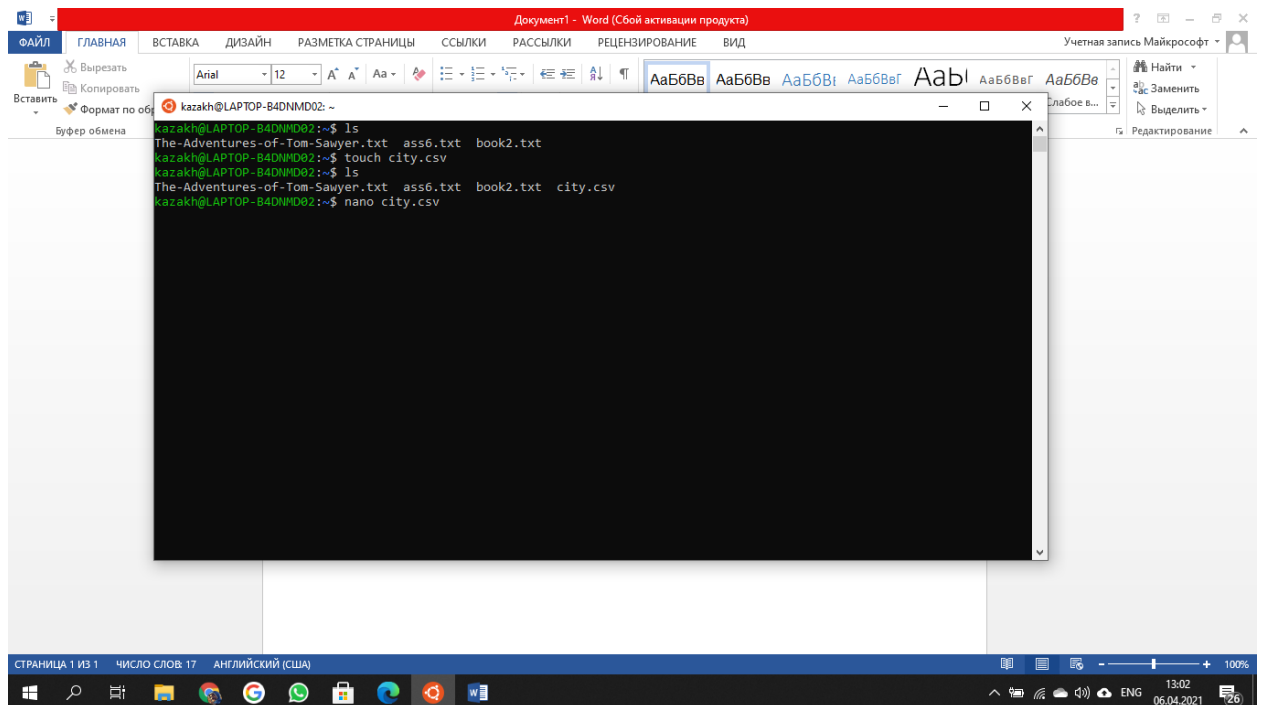
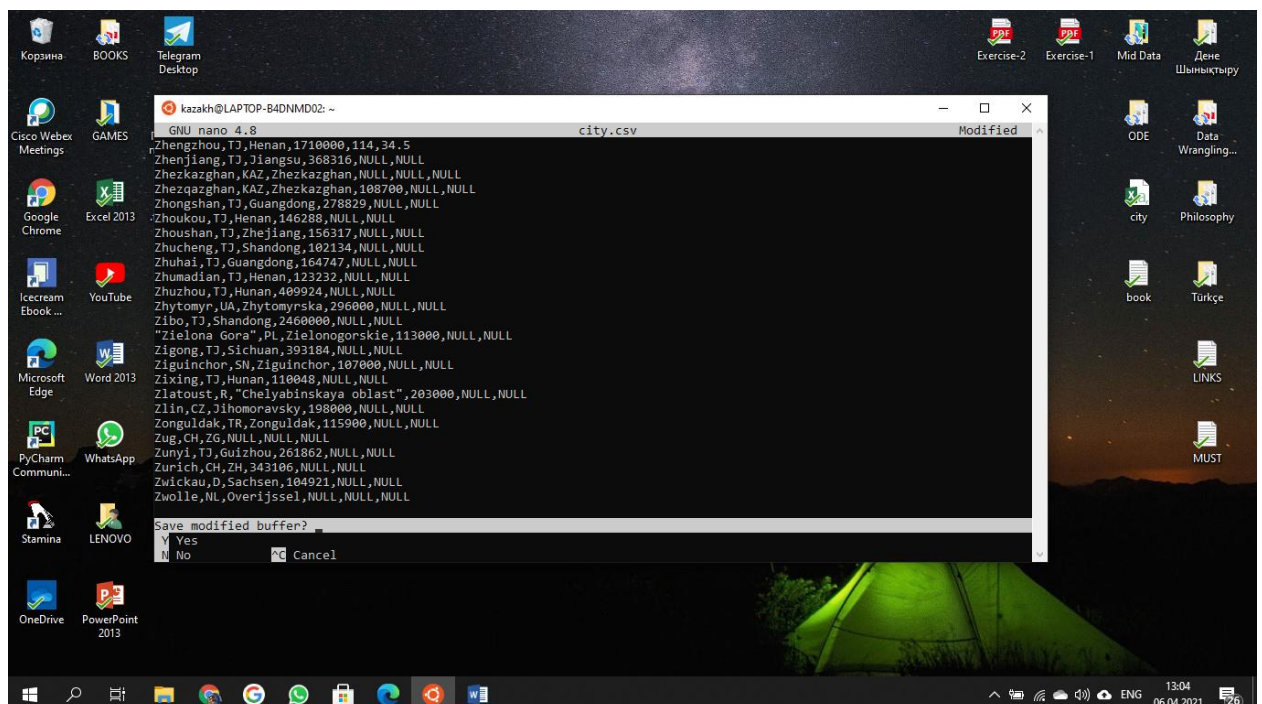# Assignment 6

## Exercise 3

**1. Create a working copy of your file city.csv (for security reasons).**

**SOLVE**

I create new csv file and named it, "city.csv".



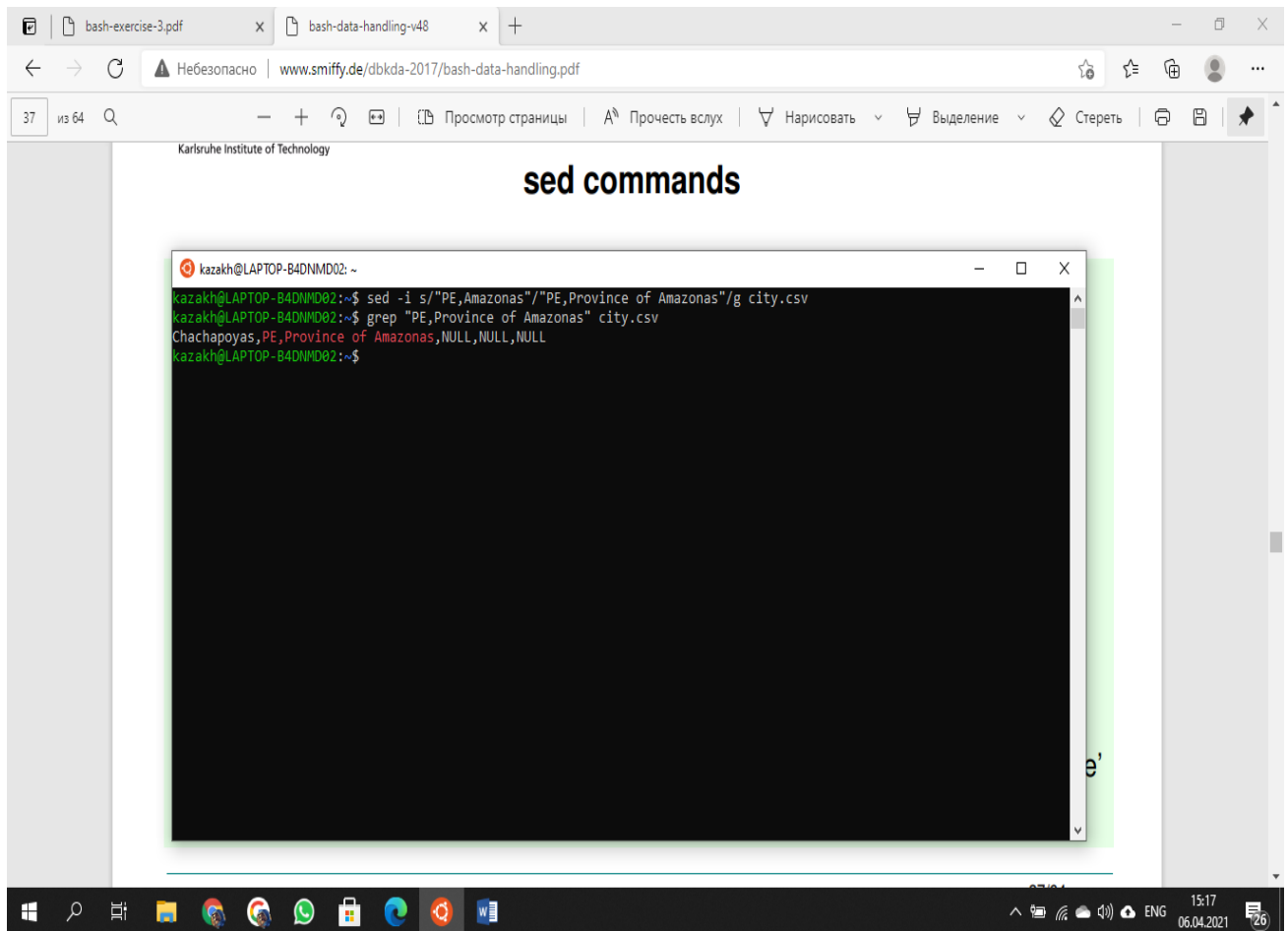And I copied all text and pasted to the csv file.

**2. Exchange in the file all occurences of the Province "Amazonas" in Peru (Code PE) with "Province of Amazonas" using sed (inplace).**

**SOLVE**

I used "sed" and –i for inplace. I replaced "Amazonas" in Peru (Code PE) with Province of Amazonas. And I checked it with grep.

**Code:** *sed –i s/"PE,Amazonas"/"PE,Province of Amazonas"/g city.csv*
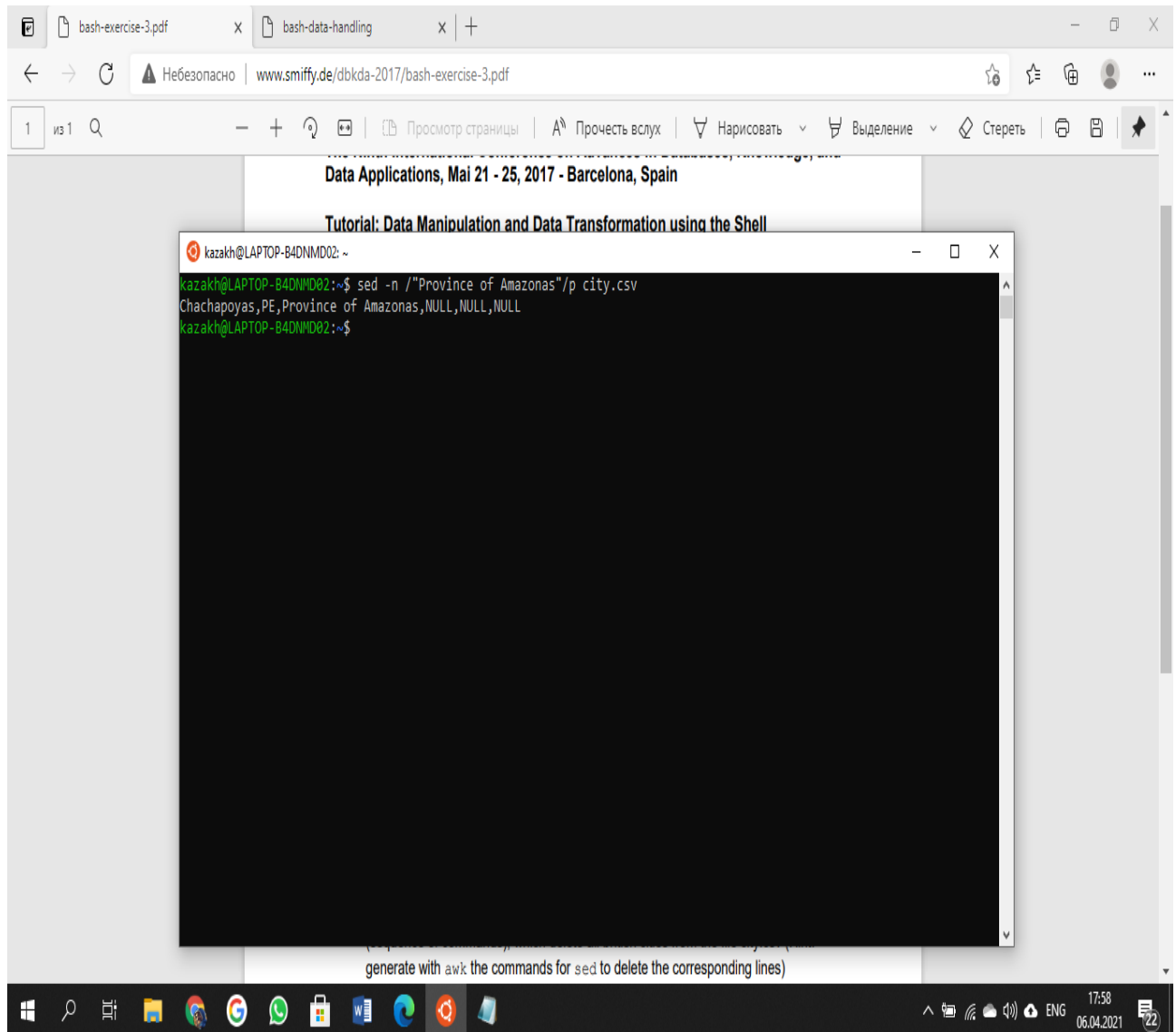
**3. Look for entries with the String "ce of Amazonas" - it should be only 1 !**

**SOLVE**

I used sed and entered *sed –n /"Province of Amazonas"/p city.csv*

It works like a grep.
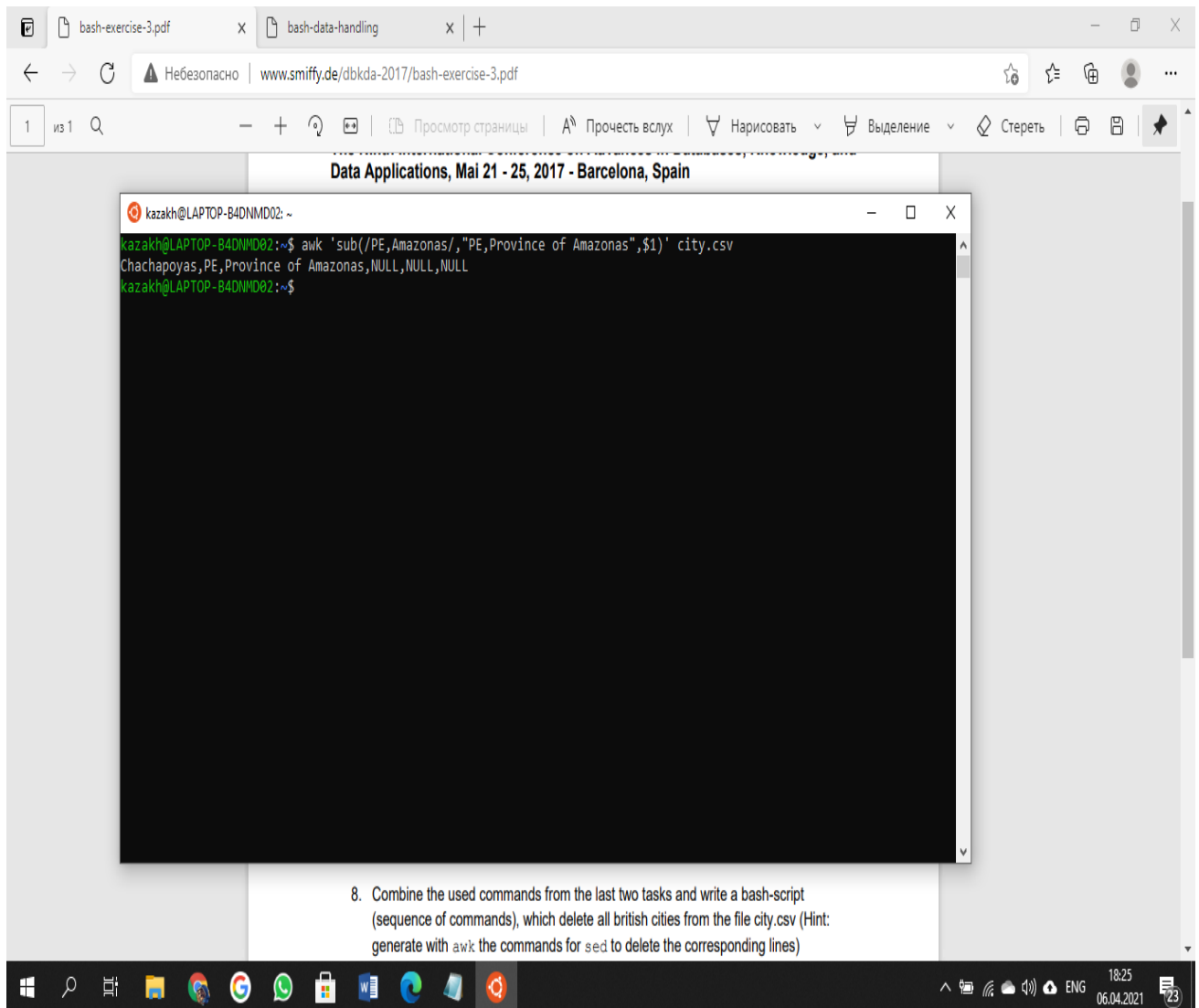
**Code:** *sed –n /"Province of Amazonas"/p city.csv*

**4. Make the same operations using awk.**

**SOLVE**

I did it with "awk"

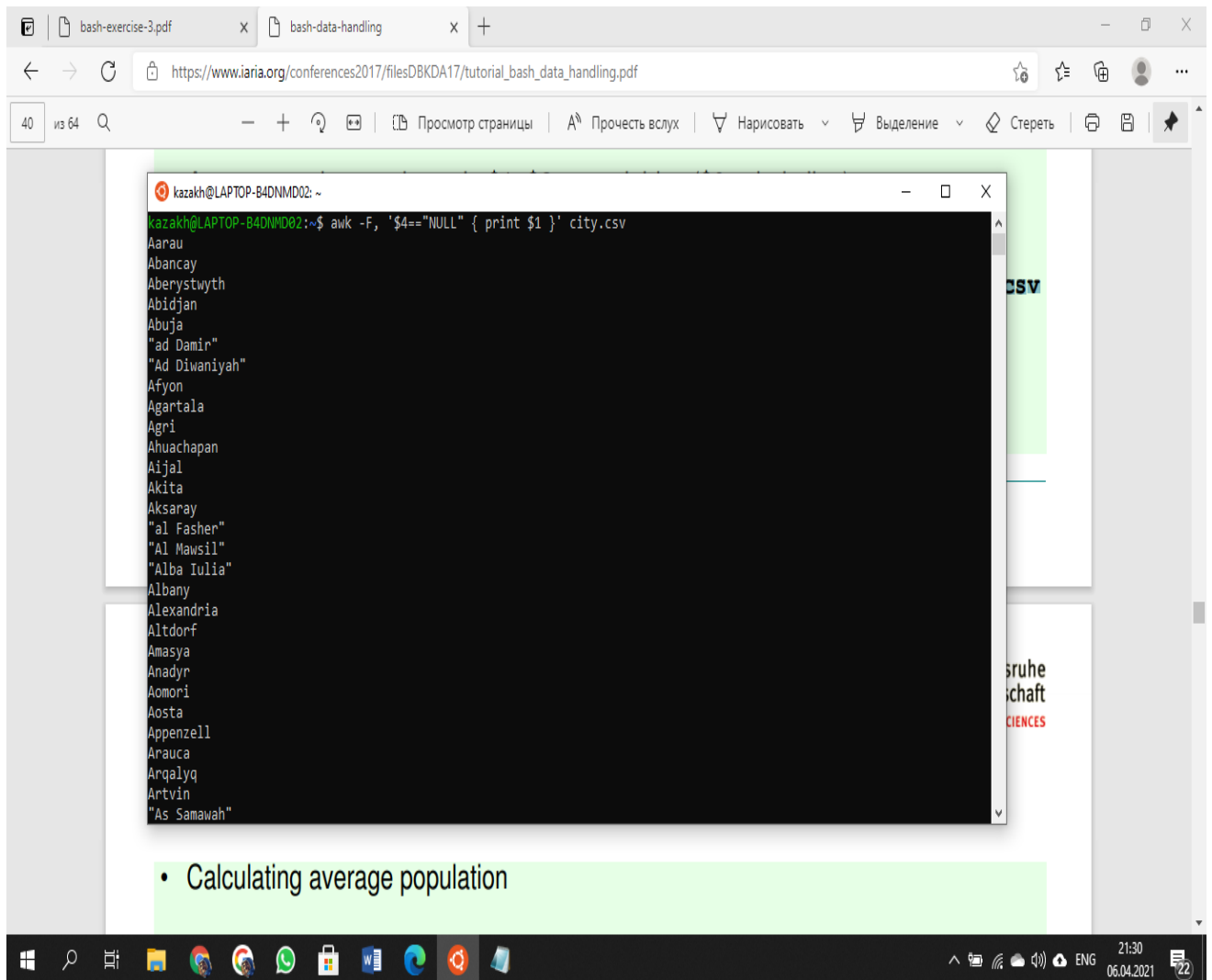**Code:** *awk 'sub(PE,Amazonas/,"PE,Province of Amazonas", $1)' city.csv*

**5. Print all cities which have no population given.**

**SOLVE**

I use awk. If fourth column equals to zero, then it print name of city.

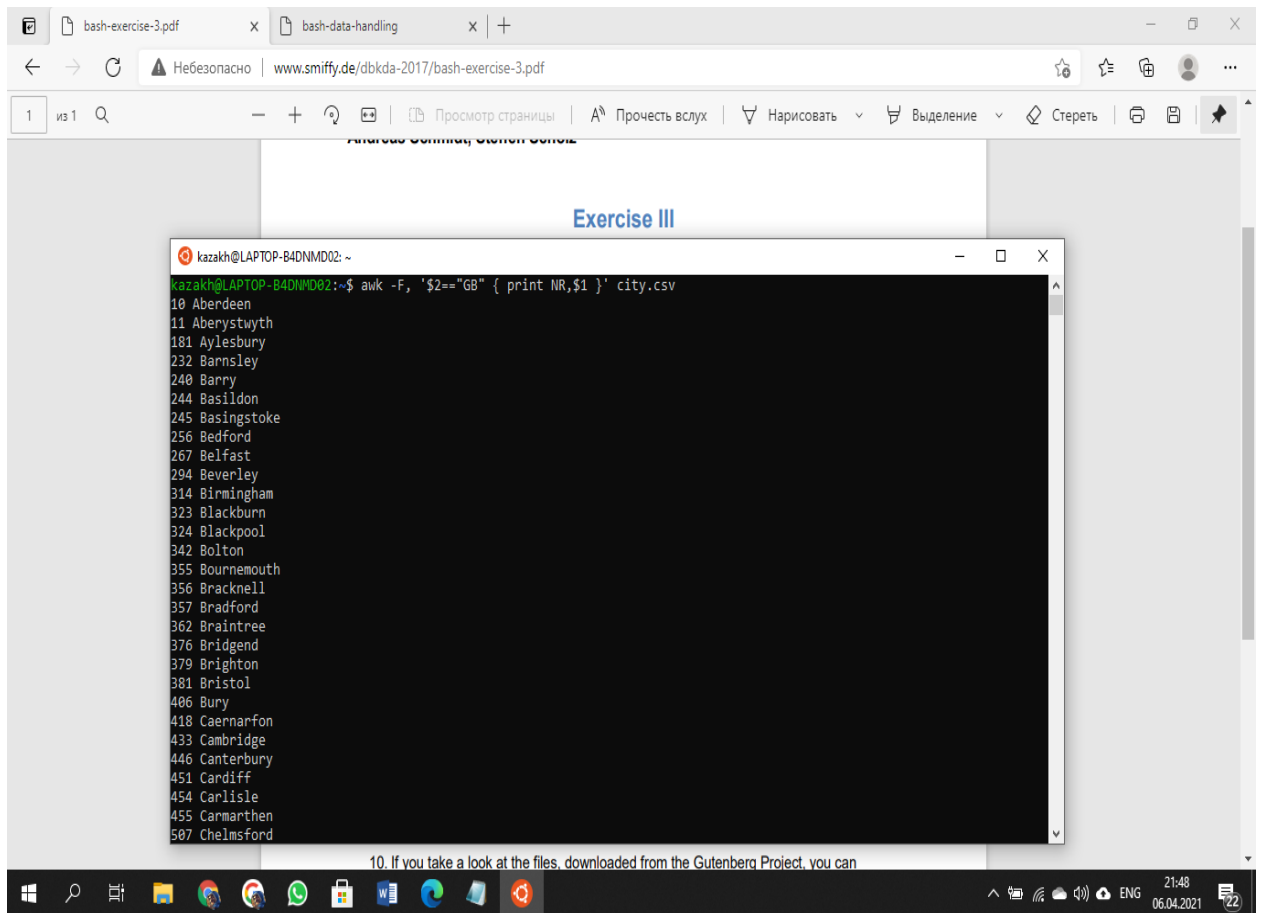**Code:** *awk –F, '$4=="NULL" { print $1 }' city.csv*

**6. Print the line numbers of the cities in Great Britain (Code: GB) using awk.**

**SOLVE**

I use awk. If second column is GB, then it print number of line and city.

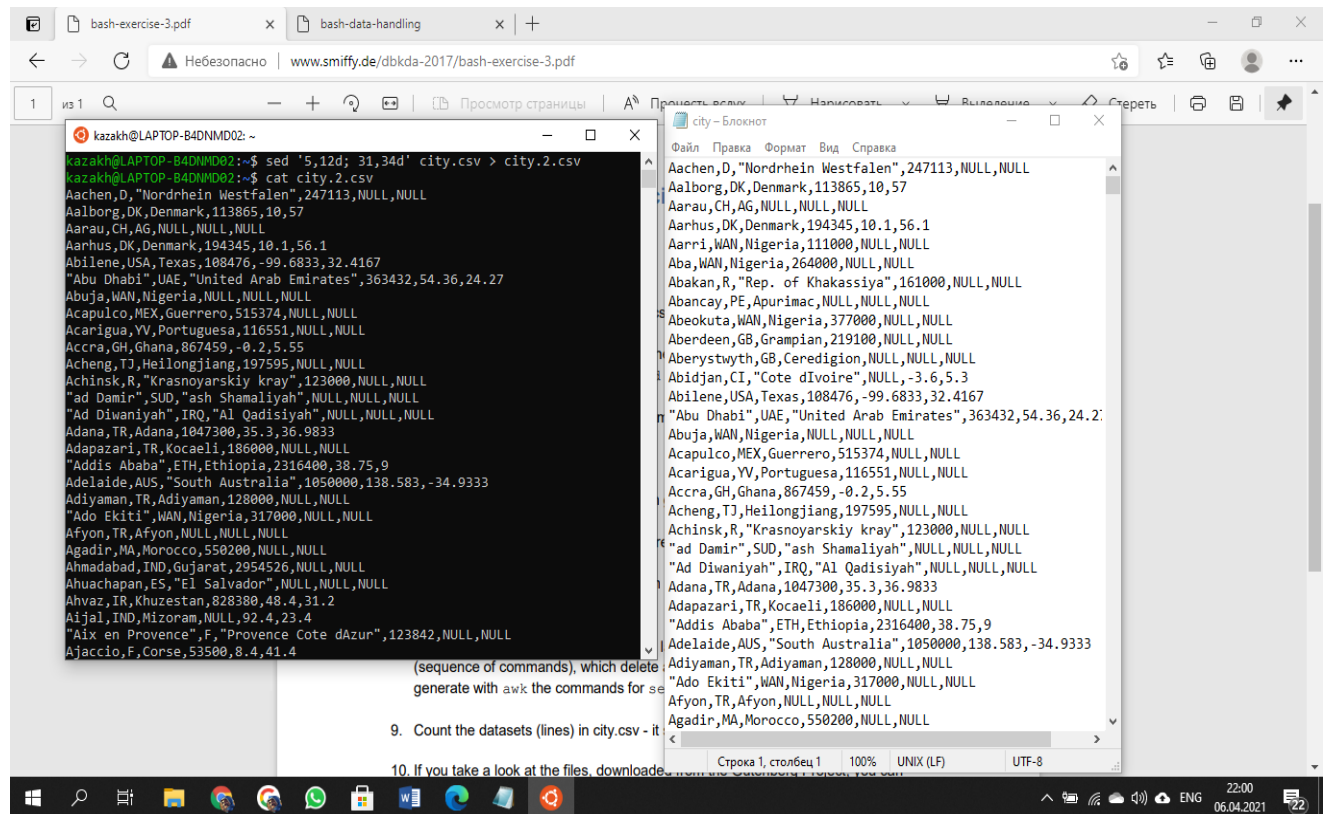**Code:** *awk  -F, '$2=="GB" { print NR,$1 }' city.csv*

## 7. Delete the records 5-12 and 31-34 from city.csv and store the result in city.2.csv using sed.

## SOLVE

I use sed. I deleted 5-12 lines and 31-34 lines. And saved result to city.2.csv.

**Code:** *sed '5,12d; 31,34d' city.csv > city.2.csv*

**8. Combine the used commands from the last two tasks and write a bash-script (sequence of commands), which delete all british cities from the file city.csv (Hint: generate with awk the commands for sed to delete the corresponding lines)**

**SOLVE**

I tried a lot, but could not find a solution. But instead of using sed and awk equally, I wrote code other ways. There are 2 ways. With awk and with sed.

**Code 1:** sed –i '/GB/ d' city.csv

**Code 2:** awk '!/'GB'/' city.csv

**9. Count the datasets (lines) in city.csv  - it should be 2880.**
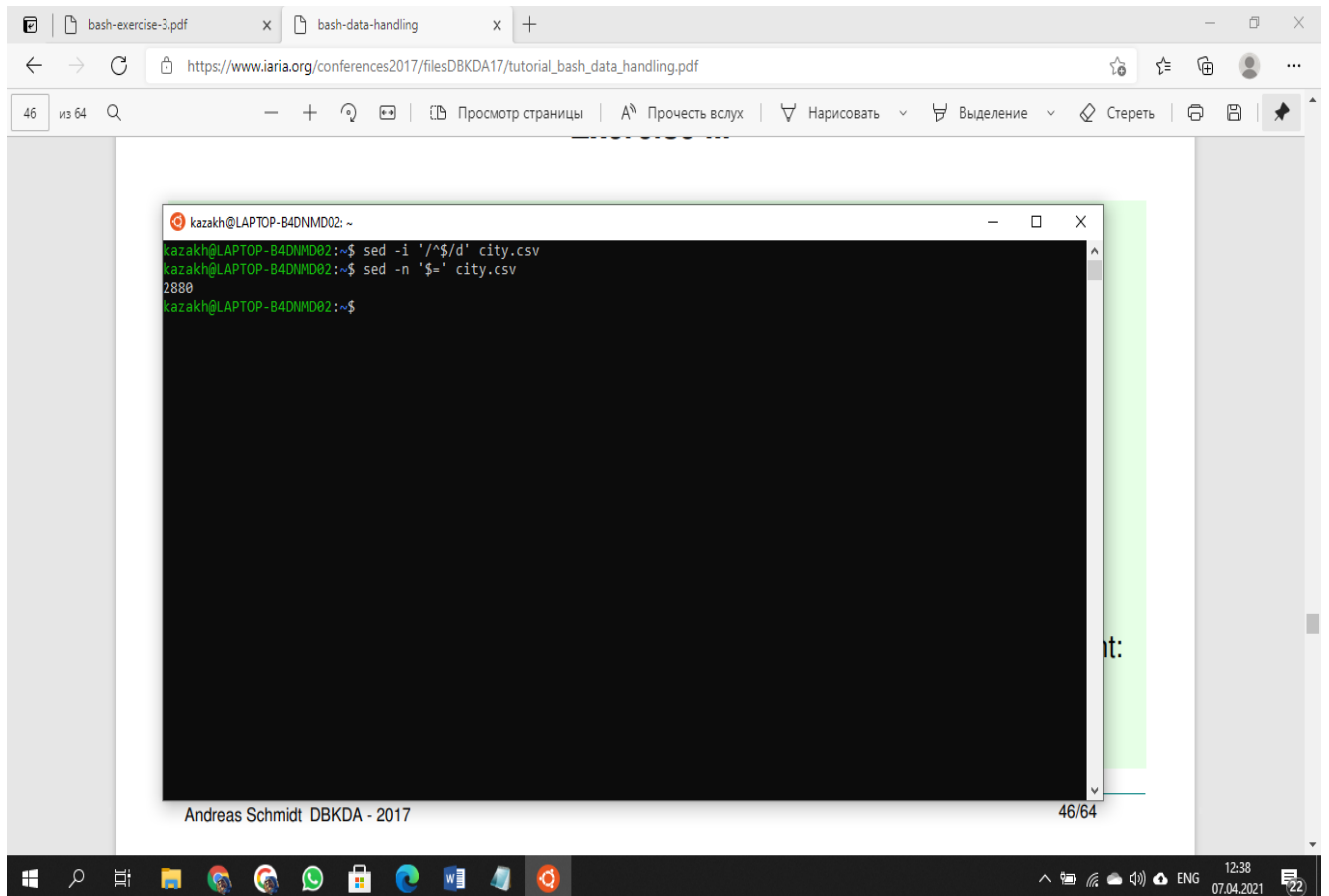

**SOLVE**

I use sed for counting. Firstly I removed all empty lines. And counted the dataset.

**Code1:** *sed –i '/^$/d' city.csv*

**Code2:** *sed –n '$=' city.csv*

**10. If you take a look at the files, downloaded from the Gutenberg Project, you can identify some boilerplate text at the begin and the end of the book. Which are the lines, who separate the literary text from the boilerplate text?**


**SOLVE**


**9. Count the datasets (lines) in city.csv - it should be 2880**


**SOLVE**