

Lab8. Text Processing Tools and Regular Expressions

Instructor :Murad Njoum



Objectives

After completing this lab, the student should be able to:

- Identify and use filters as valuable text processing tools.
- Use simple regular expressions to make text processing more efficient



Text Processing using Filters

In the pipes lab, we mentioned a group of commands called filters. These are basically commands that take some input and then filter it to produce the requested output without changing the original source of input. In this lab we will practice how to use some filters as useful tools for **text processing**

Filters:



head and tail: used to display lines from the beginning or end of a given input respectively.

cat: used to view or concatenate files.

grep: used to extract certain rows (lines) from a given input. We will concentrate on the options -i, -l (EL), -v.

cut: used to extract certain columns from a given input. We will use the options -d, -f, and -c.

tr: translates (changes) a given input to a specified output

wc: used to count lines, words, or characters in a given input.

sort: used for sorting a given input. We will present the options -i, -o, -u, -n, -k, and -t.

sed: used for stream editing (changing parts of an input to a specified output)

Create Students file(using vi)

ah6:506:Ahmad_Hamdan
sh5:345:Suha_HAMDAN
rd7:427:Ribhi_ahmad
hr4:234:hamdan_ribhi
ad6:386:Arwa_Ahmad
ad5:285:ahmadi_Ahmad



Execute the following commands:

head -2 students



Try

head -n 2 students

Or head -n +2 students

Try

head -n -2 students

```
mnjourn@ubuntu:~$ head -2 students
```

```
ah6:506:Ahmad_Hamdan
```

```
sh5:345:Suha_HAMDAN
```

View first 2 lines

```
mnjourn@ubuntu:~$ head -n 2 students
```

```
ah6:506:Ahmad_Hamdan
```

```
sh5:345:Suha_HAMDAN
```

```
mnjourn@ubuntu:~$ head -n -2 students
```

```
ah6:506:Ahmad_Hamdan
```

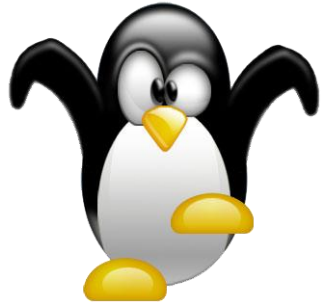
```
sh5:345:Suha_HAMDAN
```

```
rd7:427:Ribhi_ahmad
```

```
hr4:234:hamdan_ribhi
```

Execute the following commands:

- ***tail -3 students***



Try

tail -n 3 students

Try

tail -n -3 students

Try

tail -n +3 students

```
mnjourn@ubuntu:~$ tail -3 students
```

```
hr4:234:hamdan_ribhi
```

```
ad6:386:Arwa_Ahmad
```

```
ad5:285:ahmadi_Ahmad
```

View last 3 lines from file

```
mnjourn@ubuntu:~$ tail -n 3 students
```

```
hr4:234:hamdan_ribhi
```

```
ad6:386:Arwa_Ahmad
```

```
ad5:285:ahmadi_Ahmad
```

```
mnjourn@ubuntu:~$ tail -n +3 students
```

```
rd7:427:Ribhi_ahmad
```

```
hr4:234:hamdan_ribhi
```

```
ad6:386:Arwa_Ahmad
```

```
ad5:285:ahmadi_Ahmad
```

Execute the following commands:

What command would you use to get the fourth line only from file students ?

```
mnjourn@ubuntu:~$ head -4 students | tail -1  
hr4:234:hamdan_ribhi
```

```
mnjourn@ubuntu:~$ head -n 4 students | tail -n 1  
hr4:234:hamdan_ribhi
```

```
mnjourn@ubuntu:~$ head -n +4 students | tail -n -1  
hr4:234:hamdan_ribhi
```

Note: default for value on n is 10



Execute the following commands:

- *cat students*

grep ahmad students

Join both cat and grep with pipes to get the same result as the previous grep command:



`$ cat students`

```
ah6:506:Ahmad_Hamdan
sh5:345:Suha_HAMDAN
rd7:427:Ribhi_ahmad
hr4:234:hamdan_ribhi
ad6:386:Arwa_Ahmad
ad5:285:ahmadi_Ahmad
```

`$ grep ahmad`

```
students
rd7:427:Ribhi_ahmad
ad5:285:ahmadi_Ahmad
```

`$cat students | grep ahmad`

```
rd7:427:Ribhi_ahmad
ad5:285:ahmadi_Ahmad
```


grep -i Ahmad students

```
mnjoum@ubuntu:~$ grep -i ahmad students
ah6:506:Ahmad_Hamdan
rd7:427:Ribhi_ahmad
ad6:386:Arwa_Ahmad
ad5:285:ahmadi_Ahmad
```

-i: ignore-case

grep -l Ribhi * (* :means all files in current directory)

students

-l, --files-with-matches

Describe Output? Give a solution for this case ? Try *grep -L Ribhi **



***It's time to stop
posting anime***

- ***grep -v Ribhi students***

```
ah6:506:Ahmad_Hamdan  
sh5:345:Suha_HAMDAN  
hr4:234:hamdan_ribhi  
ad6:386:Arwa_Ahmad  
ad5:285:ahmadi_Ahmad
```

-v, --invert-match

Invert the sense of matching, to select non-matching lines.

grep -iv hamdan students

```
rd7:427:Ribhi_ahmad  
ad6:386:Arwa_Ahmad  
ad5:285:ahmadi_Ahmad
```



Execute the following commands:

- *cut -d: -f2 students*

```
cut -d: -f2 students
```

```
506
```

```
345
```

```
427
```

```
234
```

```
386
```

```
285
```



Q: What command would you use to get the last names for all users in file students:

```
cut -d : -f3 students | cut -d _ -f2
```

```
Hamdan  
HAMDAN  
ahmad  
ribhi  
Ahmad  
Ahmad
```



Q: What command would you use to get the first names of all users with last name **hamdan** (all cases)

```
cut -d : -f3 students | grep -i _hamdan | cut -d _ -f1
```

```
Ahmad  
Suha
```

STUDENTS-HUB.com

```
ah6:506:Ahmad_Hamdan  
sh5:345:Suha_HAMDAN  
rd7:427:Ribhi_ahmad  
hr4:234:hamdan_ribhi  
ad6:386:Arwa_Ahmad  
ad5:285:ahmadi_Ahmad
```

Execute the following commands:

cut -c2,3 students

```
cut -c2,3 students
```

```
h6
```

```
h5
```

```
d7
```

```
r4
```

```
d6
```

```
d5
```

-c, --characters=LIST
select only these characters



Continue

- *What command would you use to get the middle digit in the id numbers for all users with last name hamdan ?*

```
grep hamdan students | cut -d : -f2 | cut -c2
```

- *tr "a-z" "A-Z" < students (Describe output)*

- *What command would you use to get the first names (all in lower case) of all users that have the word **ahmad** (all cases) as part of their full name:*

```
grep -i ahmad students | cut -d : -f3 | cut -d _ -f1 | tr A-Z a-z
```

- *wc -l students*

```
wc -l students  
6 students
```



- **head -1 students | cut -d: -f3 | cut -d_ -f2 | wc -c**

7 count number of characters in last name of first student in the file students

same as previous, try with wc -w ,what does mean -w,-l,-c,-m

- **What command would you use to count the number of files in your home directory?**

ls | wc -l

sort students (Describe output)

alphabetical order

sort -o result students (What happened?)

alphabetical order save to result file

Try : sort -r result students or sort -r students > result



- *sort -k2 -t: -n students (Describe output)*

sort -k2 -t: -n students

hr4:234:hamdan_ribhi
ad5:285:ahmadi_Ahmad
sh5:345:Suha_HAMDAN
ad6:386:Arwa_Ahmad
rd7:427:Ribhi_ahmad
ah6:506:Ahmad_Hamdan

sort -k2 -t: -n students

-k2 : according to key 2 : Colum 2

-t: sperator (:)

-n: compare according to string numerical value



- What command would you use to list all the last names of users in file students sorted based on lower case letters and **without repetition**

```
cut -d : -f3 students | cut -d _ -f2 | sort -f -u
```

-f : fold lower case to upper case characters

-u :unique

```
sed 's/ahmad/damha/' students
```

```
sed s/ahmad/damha/ students
```

```
ah6:506:Ahmad_Hamdan
```

```
sh5:345:Suha_HAMDAN
```

```
rd7:427:Ribhi_damha
```

```
hr4:234:hamdan_ribhi
```

```
ad6:386:Arwa_Ahmad
```

```
ad5:285:damhai_Ahmad
```



- **What is different when we run the same command with the *i* (ignore case) option, as follows:**
sed 's/ahmad/damha/i' students

```
sed s/ahmad/damha/i students
```

```
ah6:506:damha_Hamdan
```

```
sh5:345:Suha_HAMDAN
```

```
rd7:427:Ribhi_damha
```

```
hr4:234:hamdan_ribhi
```

```
ad6:386:Arwa_damha
```

```
ad5:285:damhai_Ahmad
```



What is different when we run the same command with the g (global) option, as follows:

***sed 's/ahmad/damha/ig' students** ,*

```
sed s/ahmad/damha/ig students  
ah6:506:damha_Hamdan  
sh5:345:Suha_HAMDAN  
rd7:427:Ribhi_damha  
hr4:234:hamdan_ribhi  
ad6:386:Arwa_damha  
ad5:285:damhai_damha
```



Here the "s" specifies the substitution operation. The "/" are delimiters.

SED command in UNIX is stands for stream editor and it can perform lot's of function on file like, searching, find and replace, insertion or deletion.

Though most common use of SED command in UNIX is for substitution or for find and replace. By using SED you can edit files even without opening it, which is much quicker way to find and replace something in file, than first opening that file in VI Editor and then changing it.

- SED is a powerful text stream editor. Can do insertion, deletion, search and replace(substitution).

- SED command in Unix supports regular expression which allows it perform complex pattern matching

Regular Expressions

- Some of the filters mentioned above such as **grep** and **sed** may use what we call regular expressions to be more powerful and precise. To get more information about the power and extent of regular expressions, you can read the man pages using the command:
man regex
- **pattern\$** : applied to a pattern if it is at the **end of a given line.**
^pattern: applied to a pattern if it is at the beginning of a given line.
[abc]: means a or b or c
[^abc]: means all characters except a, b, or c.



Command Cont...

- `grep -i 'hamdan$' students`
- `cut -d: -f3 students | grep -i '^ahmad'`
- `cut -d: -f3 students | cut -d_ -f1 | grep -i '^ahmad$'`
- `cut -d: -f1 students | grep a[dh][^6]`
- `cut -d: -f3 students | sed 's/^ahmad/sameer/ig'`
- `sed 's/ahmad$/Sameer/i' students`





Puzzle Quiz Game in class

Directions:

1. Quiz is practical (at your machine in lab).
2. It's open book or notes, internet not allowed.
3. True run commands are only accepted.
4. Points are:
 1. 1 points
 2. 2 points
 3. 3 points
6. Time expired within 3 minutes, not extension allowed.
7. 1st, 2nd students whom complete the task will get full mark, others will loose marks (-1,-2,-3,...etc.)

(2 points) Display the default shell used by user root

```
grep ^root /etc/passwd | cut -d : -f7 | cut -d : -f2
```

(2 points) Display the number of files in directory /etc that end with the word .conf

```
ls /etc | grep .conf$ | wc -l
```

(1 points) Print lines form 15-17 from /etc/passwd file?

```
head -n 17 /etc/passwd | tail -n 3
```

(2 points) List the login names for all users with the bash as their default shell.

```
grep bash$ /etc/passwd | cut -d : -f1
```

Questions: (use passwd file)

(**3 points**) Display the first names of all users whose last names end with the letter 'n' or 'm' (for all cases):

```
grep ^u1 /etc/passwd | cut -d : -f5 | grep -i [mn]$
```

(**3 points**) Display the last names of all the users sorted by their user id numbers (**descending order**)

```
grep ^u1 /etc/passwd | sort -r -k3 -t: -n | cut -d : -f5 | cut -d _ -f2
```

(**3 points**) Display all files in (etc) directory does not contain word passwd in these files page by page at screen

```
grep -L passwd /etc/* 2>outerror | more
```






exciting
EDUCATION
INFORM ENRICH INSPIRE