Défi BigData Information Systems for Big Data 1 - Written Exam -

Luis Gustavo Nardin <gnardin@emse.fr>

18 April 2025

Duration: 2 hours (9h00-11h00)

Instructions

- You are allowed to use an A4 sheet of paper, handwritten by you.
- You are **forbidden** to use calculators, telephones, computers or any other means of communication.
- Write your **full name** in the first page of this document.
- Write your **answers directly on this document** and hand it in in its entirety at the end of the allotted time.
- The multiple choice questions have only one correct answer.

Name			
Grade:			
Remarks:			

${\bf Question}~{\bf 1}:~(1~{\rm point})$ What is who in the command LANG=C man -k who:
\square an argument to the option $-k$
\square a long option
\Box an argument to the program
\square a relative path
Question 2: (1 point) When executing a command, the command is searched in the directories listed in the PATH environment variable :
\Box from left to right
\square in no particular order
\Box from right to left
\Box including only the first 256 characters
Question 3: (1 point) Which is NOT a valid command to declare the environment variable ENV_VAR?
□ export ENV_VAR=1
□ newvar -x ENV_VAR=1
☐ ENV_VAR=1
☐ declare -x ENV_VAR=1
Question 4: (1 point) Which command allows us to copy the file notes.txt located in the directory tmp in the current directory to the directory /home/user with the name new_notes.txt?
☐ cp /tmp/notes.txt /home/user/new_notes.txt
☐ cp /tmp/notes.txt home/user/new_notes.txt
☐ cp tmp/notes.txt /home/user/new_notes.txt
☐ cp tmp/notes.txt /new_notes.txt
Question 5: (1 point) The #!/bin/bash in the first line of a Shell script
\square specifies the path of the file in which the result of the script will be written
\square specifies the path where the script is stored
\square is a comment
\square specifies the executable shell to use for running the script

Question 6: (1.5) Write the regular expression that matches all the patterns below.

Hello World hello world world hello WORLD HELLO hello hello world hello world world

Question 7: (1.5 points) The file games.txt in your current directory has the following content

2024-08-16	Nimes vs Orléans	Loss	1-2
2024-08-23	Nimes vs Châteauroux	Win	1-0
2024-08-30	Valenciennes vs Nimes	Win	2-0
2024-09-06	Nimes vs Rouen	Match	2-2
2024-09-13	Concameau vs Nimes	Win	1-0
2004-09-20	Nimes vs Versailles	Win	1-0

Write a command or chain of commands (NOT a script) to find all games in which *Nimes* won when playing at home. The team playing at home is the team name before "vs".

Question 8: (2 points) The command who returns the list of users logged on in the computer in the format

Write a command or chain of commands (NOT a script) to create the file logged_on.txt which contains a listing of the unique user names currently logged on in inverted alphabetical order.

Question 9: (3.0 points) Suppose you have the users.csv file with the content below. id,first-name,last-name,email,login,logout,cost

- 1, John, Smith, john.smith@example.com, 2023-01-15, 2023-03-20, 125.99
- 2, Jane, Doe, jane. doe@example.com, 2023-01-16, 2023-03-21, 210.50
- 3,Bob,Johnson,bob@example.com,2023-01-17,2023-03-22,0
- 4, Alice, Williams, alice. williams@example.com, 2023-01-18, 75.25
- 5, Brown, mike. brown@example.com, 2023-01-19, 2023-03-24, 150.75
- 6,Sarah,Miller,sarah.miller@example.com,invalid date,2023-03-25,95.00
- 7, David, Jones, david.jones@example.com, 2023-01-21, 2023-03-26, 300.00
- 8, Lisa, Garcia, lisa.garcia@example.com, 2023-01-22, 2023-03-27, -50.00
- 9, James, Martinez, mymail@example.com, 2023-01-23, 2023-03-28, 125.00

Write an awk script that

- 1. Drop the heading line (i.e., first line)
- 2. Fill in empty name fields (i.e., first-name and last-name) with the word Unknown
- 3. Check if the dates are properly formatted, replacing those not properly formatted by: login field with 2023-01-20 and logout with 2023-01-23
- 4. Replace negative costs by 0.0

Question 10: (3 points) Explain each line of code in the script below (arguments, processing, results, etc.).

```
#!/bin/bash
if [[ $# -ne 1 ]]
then
  echo "Usage: ${0} nameList" >&2
  exit 1
fi
if [[ ! -f "$1" ]]
then
 echo "${1} error." >&2
 exit 1
fi
usernames='cat "$1"'
for i in $usernames
do
  if [[ -d $i ]]
    echo "${i} already exists." >&2
    exit 1
  fi
  mkdir $i
  chmod 775 $i
  chown $i $i
done
```

Question 11: (4 points) You have been hired to provide a solution to summarize and consolidate the data of the daily sales generated by the stores of a large retailer. The stores are geographically distributed all over the world and the headquarter is located in Saint-Étienne. The stores generate one file per day containing their sales in the format below.

```
\label{lem:price} \begin{split} &\text{Date\_Time;Store;Product;Quantity;Unit\_Price;Unit\_Cost;Unit\_Profit;Total\_Price} \\ &02/10/2022 \ 08:01; store10; p3;7;269.33;105.20;164.13;1885.28 \\ &02/10/2022 \ 08:49; store10; p32;9;38.32;15.21;23.10;344.85 \\ &02/10/2022 \ 08:54; store10; p2;4;414.34;211.85;202.49;1657.35 \\ &\dots \end{split}
```

You are required to write a Bash script (no awk script only) that receives as input the directory where the daily files are stored. The script should check if the directory exists and return an error if not. Moreover, the script must contain a function that receives a full path to a file containing the sale records, summarizes the daily sales by product and writes the result to the same directory with name summary.txt.

Command References

OPTIONS

CAT(1) User Commands CAT(1) Pattern Syntax -E, --extended-regexp NAME Interpret PATTERNS as extended regular expressions (EREs, see below). cat - concatenate files and print on the standard output SYNOPSIS Matching Control cat [OPTION]... [FILE]... -v, --invert-match Invert the sense of matching, to select non-matching GREP(1) User Commands GREP(1) lines. NAME General Output Control grep, egrep, fgrep, rgrep - print lines that match patterns -c, --count Suppress normal output; instead print a count of SYNOPSIS matching lines for each input file. With the -v, grep [OPTION...] PATTERNS [FILE...] --invert-match option (see above), count non-matching grep [OPTION...] -e PATTERNS ... [FILE...] lines. DESCRIPTION -o, --only-matching grep searches for PATTERNS in each FILE. PATTERNS is one or Print only the matched (non-empty) parts of a matching more patterns separated by newline characters, and grep prints line, with each such part on a separate output line. each line that matches a pattern. Typically PATTERNS should be quoted when grep is used in a shell command. SORT(1) User Commands SORT(1) A FILE of "-" stands for standard input. If no FILE is given, recursive searches examine the working directory, and NAME nonrecursive searches read standard input. sort - sort lines of text files Debian also includes the variant programs egrep, fgrep and SYNOPSIS rgrep. These programs are the same as grep -E, grep -F, and sort [OPTION]... [FILE]... sort [OPTION]... --files0-from=F grep -r, respectively. These variants are deprecated upstream, but Debian provides for backward compatibility. For portability reasons, it is recommended to avoid the variant programs, and DESCRIPTION use grep with the related option instead. Write sorted concatenation of all FILE(s) to standard output.

With no FILE, or when FILE is -, read standard input.

 ∞

NAME

SYNOPSTS

tr - translate or delete characters

Mandatory arguments to long options are mandatory for short options too. Ordering options:

-d, --dictionary-order consider only blanks and alphanumeric characters

-h, --human-numeric-sort compare human readable numbers (e.g., 2K 1G)

-r, --reverse reverse the result of comparisons

-t, --field-separator=SEP use SEP instead of non-blank to blank transition

-u, --unique
 with -c, check for strict ordering; without -c, output
 only the first of an equal run

TEE(1) User Commands TEE(1)

tee - read from standard input and write to standard output and files

tee [OPTION]... [FILE]...

TR(1) User Commands TR(1)

NAME

SYNOPSIS

tr [OPTION]... STRING1 [STRING2]

DESCRIPTION

Translate, squeeze, and/or delete characters from standard input, writing to standard output. STRING1 and STRING2 specify arrays of characters ARRAY1 and ARRAY2 that control the action.

-c, -C, --complement use the complement of ARRAY1

-d, --delete
 delete characters in ARRAY1, do not translate

-s, --squeeze-repeats replace each sequence of a repeated character that is listed in the last specified ARRAY, with a single occurrence of that character

-t, --truncate-set1 first truncate ARRAY1 to length of ARRAY2

--help display this help and exit

--version

output version information and exit

ARRAYs are specified as strings of characters. Most represent themselves. Interpreted sequences are:

\NNN character with octal value NNN (1 to 3 octal digits)

\\ backslash

\a audible BEL

\b backspace

f form feed

NAME

rence.

writing to OUTPUT (or standard output).

With no options, matching lines are merged to the first occur-

\n new line -c, --count \r return prefix lines by the number of occurrences \t horizontal tab \v vertical tab Note: 'uniq' does not detect repeated lines unless they are adjacent. You may want to sort the input first, or use 'sort -u' CHAR1-CHAR2 without 'uniq'. all characters from CHAR1 to CHAR2 in ascending order [CHAR*] in ARRAY2, copies of CHAR until length of ARRAY1 WC(1) User Commands WC(1) [CHAR*REPEAT] REPEAT copies of CHAR, REPEAT octal if starting with 0 NAME Γ=CHAR=] wc - print newline, word, and byte counts for each file all characters which are equivalent to CHAR SYNOPSIS Translation occurs if -d is not given and both STRING1 and wc [OPTION] ... [FILE] ... STRING2 appear. -t may be used only when translating. ARRAY2 wc [OPTION]... --filesO-from=F is extended to length of ARRAY1 by repeating its last character as necessary. Excess characters of ARRAY2 are ignored. Char- DESCRIPTION acter classes expand in unspecified order; while translating, Print newline, word, and byte counts for each FILE, and a total [:lower:] and [:upper:] may be used in pairs to specify case line if more than one FILE is specified. A word is a conversion. Squeezing occurs after translation or deletion. non-zero-length sequence of printable characters delimited by white space. UNIQ(1) User Commands UNIQ(1) With no FILE, or when FILE is -, read standard input. The options below may be used to select which counts are printed, always in the following order: newline, word, characuniq - report or omit repeated lines ter, byte, maximum line length. SYNOPSIS uniq [OPTION]... [INPUT [OUTPUT]] -c, --bytes print the byte counts DESCRIPTION Filter adjacent matching lines from INPUT (or standard input), -m, --chars

print the character counts

print the newline counts

-l. --lines