

COMP1204

Unix Coursework

1 OVERVIEW

This coursework will cover two key topics that will have been covered in the first three weeks of the module: Unix and Latex. The coursework is divided into two parts: (i) Unix scripting for basic file processing and data analysis and (ii) Report writing using Latex. Each part carries a percentage of the marks for this coursework (out of a 100) as detailed in Table 1. You should be able to complete various parts of this coursework as we go through the lectures, and parts of this assignment will be covered during lab sessions. Help will also be available throughout the labs. The key points are:

- This Coursework counts for 10% of this module.
- The deadline¹ for submission of your scripts and report: **12th March 2021 by 4pm.**
- Feedback will be given within 4 weeks after the deadline.
- You are only allowed to use Bash (Unix) commands and Latex (Overleaf and other editors are acceptable). Use of other scripting languages or text editors will get you zero marks for the relevant sections.

For submission instructions, please see Section 4 at the end of this document.

¹late penalty.

2 LEARNING OUTCOMES (LOS)

This problem sheet aims to achieve the following learning outcomes:

- Knowledge of Unix commands and pipes and filters.
- Knowledge of Latex commands and document preparation.
- Data cleaning techniques using pattern matching and filtering

3 THE ASSIGNMENT

You work as a data scientist for TripAdvisor and your job is to help them make sense of what hotels are performing well. You have been tasked with analysing all the files containing reviews for each hotel as described in the next subsections.

Before you start, create a project on <https://git.soton.ac.uk>. The project should be named `comp1204-userid-cwk1`. So if your userid is `je5g15`, the project should be named `comp1204-je5g15-cwk1`. **Make sure the visibility level of the project is set to private. Now invite user je5g15 as Maintainer on the project.**

Please make sure to regularly commit your work on the Unix scripts to your git project.

3.1 DATASET

The dataset to be used for this coursework is the TripAdvisor dataset at: https://secure.ecs.soton.ac.uk/notes/comp1204/coursework/dataset/reviews_dataset.tar.gz. Download this file to a folder on your home drive (e.g., `myworkspace`). Extract the file using standard UNIX file decompression commands.

3.1.1 BASIC FILE PROCESSING AND DATA ANALYSIS – 90%

You would like to find out what hotels are commented on the most (the frequency of comments may indicate how many guests they actually receive). To answer this question:

1. Copy a hotel data file to your home directory (e.g., `hotel_72572.dat`).
2. Write a Unix script called *countreviews.sh* that counts the number of reviews in the file that takes input from the command line like this:

```
% ./countreviews.sh hotel_72572.dat
15
```

where `hotel_72572.dat` is an example file name. Note that 15 is just an example and not the actual result for this file. Also note that the argument to the script, `hotel_72572.dat`, is just an example file PATH; your script must be able to handle both relative and absolute paths correctly. **Do NOT submit this script on the ECS electronic hand-in system – complete the next steps and submit the final script for this section.**

3. Extend *countreviews.sh* to count the number of reviews in each file given the folder name (i.e., where all your files are stored).

```
% ./countreviews.sh path_to_reviews_folder
15
12
13
...
```

Remember again that the argument to the script, *path_to_reviews_folder*, can be an absolute or a relative path; your script must be able to handle both types of paths correctly.

4. Finally rank all the hotels according to the review count so that the hotel with the most reviews is at the top of your list.

The output of your complete script should be formatted like this:

```
% ./countreviews.sh path_to_reviews_folder
hotel_1322 50
hotel_21313 49
hotel_31331 45
...
```

The argument to the script, *path_to_reviews_folder*, can be an absolute or a relative path; your script must be able to handle both types of paths correctly.

Regarding the output of the script, please note the following:

1. The *.dat* extension is omitted from the filename.
2. The hotel name is separated from the count by a single whitespace.
3. Nothing but the hotel-count pairing should be output by the script.

Submit only your final script on the ECS electronic hand-in system.

3.1.2 REPORT – 10%

Write a report in Latex detailing the following:

1. The cover page of the report have at least your name and ID written on it as well as the title.
2. The *countreviews.sh* script you wrote. Make sure you clearly explain what the script is doing. Your script should preferably be written using the 'listings' environment from Latex.

Unix Script for basic file processing and data analysis	90%
Report Writing in Latex	10%

Table 1: The weighting given to the different parts of this coursework.

3.2 ASSESSMENT CRITERIA

Your code and your report will be evaluated as follows:

Unix

For Unix scripts, we will use an automated script checker that will pipe data to your code and check the output. Your scripts will be tested on previously unseen data. In particular, we will check that:

1. Code returns expected output.
2. Code cleanliness and efficiency (slow code will be penalised).
3. Appropriate use of git versioning system in writing your scripts. In particular, we will assess if meaningful commits and commit messages have been provided.

Indicative feedback may include: code does not work, code works partially (i.e., some functions not working), code is inefficient, code not readable, git commit messages not provided, all scripts work.

Report

For the report we will check that the report is written in LATEX and any script listed using standard LATEX environments e.g., verbatim, listings, or algorithmic.

4 SUBMISSION INSTRUCTIONS

While working on your script, `countreviews.sh`, you would have made several commits on git. For your submission generate a git log file – named `git.log`.

Submit your work using the ECS electronic hand-in system. The submission is to be made by **4pm** on the due date listed above. Please submit a single file to the ECS electronic hand-in system as detailed below:

- Your script and report must be in compressed archive named as `comp1204.tar.gz`. You should submit the script (`countreviews.sh`), the git log file (`git.log`), your report in LATEX (`report.tex` and `report.pdf`), and the LATEX log file (`report.log`) that is generated when you compile your latex file (on Overleaf you need to navigate the options to download the log file).
- Your report should be in .PDF format and be included in the archive.

Failure to follow these instructions will incur a penalty. In particular, you will lose (possibly all) marks if:

1. You use Word to create your document.
2. You compile to ZIP and change the extension of the file.
3. You submit a word doc and change the extension of the file to PDF
4. Your code runs for more than 1 min on the test dataset.