

CITS 1401 Computational Thinking with Python



Project 1, Semester 1, 2024

Submission deadline: **19th April 2024, 6:00 PM**

Total Marks: **30 (Value: 15%)**

Project description

You should construct a Python 3 program containing your solution to the following problem and submit your program electronically on Moodle. The name of the file containing your code should be your student ID e.g., 12345678.py. No other method of submission is allowed. **Please note that this is an individual project.** Your program will be automatically run on Moodle for some sample test cases provided in the project sheet if you click the "check" link. However, this **does not test** all required criteria and your submission will be **manually** tested thoroughly for grading purposes after the due date. Remember you need to submit the program as a single file and copy-paste the same program in the provided text box. You have only one attempt to submit, so do not submit until you are satisfied with your attempt. All open submissions at the time of the deadline will be automatically submitted. Once your attempt is submitted, there is no way in the system to open/reverse/modify it.

You are expected to have read and understood the University's guidelines on academic conduct. In accordance with this policy, you may discuss with other students the general principles required to understand this project, but the work you submit must be the result of your own effort. Plagiarism detection, and other systems for detecting potential malpractice, will therefore be used. Besides, if what you submit is not your own work then you will have learned little and will therefore, likely, fail the final exam.

You must submit your project before the deadline mentioned above. Following UWA policy, a late penalty of 5% will be deducted for each day i.e., 24 hours after the deadline, that the assignment is submitted. No submissions will be allowed after 7 days following the deadline except approved special consideration cases.

Project Overview

In today's digital age, social media has become an integral part of our daily lives, shaping how we connect, communicate, and consume information. With millions of users engaging across various platforms, understanding user behaviours and trends on social media has become increasingly important. The aim of this project is to analyse user demographics to better understand social media usage.

The dataset for this project comprises several key columns, including age, gender, time spent on social media (in hours), platform, interests, country, demographics, profession, income, and debt status.

You are required to write a Python 3 program that will read a CSV file. After reading the file, your program is required to complete the following tasks:

- 1) Find the list of student details (ID and income) for a specific country who are in debt (or have debt status True) and spending more than 7 hours on any social media.
- 2) Find the list of unique countries for a specific age group.
- 3) Find the age statistics for a specific age group. The age statistics include average time spent in hours, standard deviation of income, and which demography (e.g.

rural, sub_urban or urban) spent the lowest average time on social media for the specific age group.

- 4) Find the platform that has the highest number of users and calculate the correlation between the age and the income for that user base.

Requirements

- 1) **You are not allowed to import any external or internal module in python.** While use of many of these modules, e.g., `csv` or `math` is a perfectly sensible thing to do in a production setting, it takes away much of the point of different aspects of the project, which is about getting practice opening text files, processing text file data, and use of basic Python structures, in this case lists and loops.
- 2) Ensure your program does NOT call the `input()` function at any time. Calling the `input()` function will cause your program to hang, waiting for input that the automated testing system will not provide (in fact, what will happen is that if the marking program detects the call(s), it will not test your code at all which may result in zero grade).
- 3) Your program should also not call `print()` function at any time except for the case of graceful termination (if needed). If your program has encountered an error state and is exiting gracefully then your program needs to return `zero` for numerical values such as average time spent, standard deviation, `None` for non-numeric output such as demography, otherwise empty list and print an appropriate message. At no point should you print the program's outputs instead of (or in addition to) returning them or provide a printout of the program's progress in calculating such outputs.
- 4) Do not assume that the input file names will end in `.csv`. File name suffixes such as `.csv` and `.txt` are not mandatory in systems other than Microsoft Windows. Do not enforce that within your program that the file must end with a `.csv` or any other extension (or try to add an extension onto the provided csv file argument), doing so can easily lead to syntax error and losing marks.

Input

Your program must define the function `main` with the following syntax:

```
def main(csvfile, age_group, country):
```

The input arguments for this function are:

- `csvfile`: The name of the CSV file (as string) containing the record of the user details. The first row of the CSV file will contain the headings of the columns. A sample CSV file "SocialMedia.csv" is provided with the project sheet on LMS and Moodle.
- `age_group`: A list parameter [`lower`,`upper`] containing the lower and upper bound of age.
- `country`: A string parameter containing the name of a country.

Output

We expect 4 outputs in the order below.

- i. **OP1** = [list of student details]: A list containing list(s) of student details [ID and income] from the `country` provided as input who are in debt and spends more than 7 hours on any social media. The student details' are ordered ascendingly based on student IDs.
- ii. **OP2** = [list of unique countries]: a list of unique countries for users whose age falls within the lower and upper bound of input `age_group` (inclusive). The list needs to be sorted in alphabetically ascending order.
- iii. **OP3** = [average time spent, standard deviation of income, demography]: A list containing a numeric value for average time spent by users whose age falls within the lower and upper bound of input `age_group` (inclusive), standard deviation of their income, and a string for demography name that represents whether rural, sub_urban or urban users falling within the `age_group` (inclusive) spent the lowest average time on social media compared to other demographics. In case two demographics have same lowest value then sort the demographics in alphabetical order and return the first one.
- iv. **OP4** = Correlation value: A numeric value for correlation between age and income for the user base of the social media platform that has the highest number of users. If there are multiple social media platforms having same highest number of users then sort them in alphabetical order and find correlation considering the first one.

标准差

All returned numeric outputs (both in lists and individual) must contain values rounded to four decimal places (if required to be rounded off). Do not round the values during calculations. Instead, round them only at the time when you save them into the final output variables.

Examples

Download `SocialMedia.csv` file from the folder of Project 1 on LMS or Moodle. An example of how you can call your program from the Python shell and examine the results it returns, is provided below:

```
>> OP1, OP2, OP3, OP4 = main('SocialMedia.csv', [18,25], 'australia')
```

The returned output variables are:

```
>> OP1
```

```
[['11', 4708.0], ['126', 5785.0], ['184', 9266.0]]
```

```
>> OP2
```

```
['australia', 'bangladesh', 'ireland', 'new zealand', 'pakistan', 'yemen']
```

```
>> OP3
```

```
[3.5556, 112446.1548, 'rural']
```

```
>> OP4
```

```
0.4756
```

Assumptions

Your program can assume the following:

- Anything that is meant to be string (e.g., header) will be a string, and anything that is meant to be numeric will be numeric in the CSV file.
- All string data in the CSV file is case-insensitive, which means "Australia" is same as "australia". Your program needs to handle the situation to consider both to be the same. Similarly, your program needs to handle the string input parameter in the same way.
- The order of columns in each row will follow the order of the headings provided in the first row. However, rows can be in random order except the first row containing the headings.
- No data will be missing in the CSV file; however, values can be zero and must be accounted for mathematical calculations.
[In case any part of the calculation cannot be performed due to zero values or other boundary conditions, do a graceful termination by printing an error message and returning a zero value (for numbers), None for (string) or empty list depending on the expected outcome. **Your program must not crash.**]
- Number of different social medial platform and country will vary, so do not hardcode.
- The main function will always be provided with valid input parameters.
- The necessary formulae are provided at the end of this document.

Important grading instruction

Note that you have not been asked to write specific functions. The task has been left to you. However, it is essential that your program defines the top-level function `main(csvfile, age_group, country)` (hereafter referred to as "main()" in the project document to save space when writing it. Note that when `main()` is written, it still implies that it is defined with its three input arguments). The idea is that within `main()`, the program calls the other functions. Of course, these functions may then call further functions. This is important because when your code is tested on Moodle, the testing program will call your `main()` function. So if you fail to define `main()`, the testing program will not be able to test your code and your submission will be graded zero. Don't forget the submission guidelines provided at the start of this document.

Marking rubric

Your program will be marked out of 30 (later scaled to be out of 15% of the final mark). 22 out of 30 marks will be awarded automatically based on how well your program completes a number of tests, reflecting normal use of the program, and also how the program handles various states including, but not limited to, different numbers of rows in the input file and / or any error or corner states/cases. You need to think creatively what your program may face. Your submission will be graded by data files other than the provided data file. Therefore, you need to be creative to investigate corner or worst cases. I have provided few guidelines from ACS Accreditation manual at the end of the project sheet which will help you to understand the expectations.

8 out of 30 marks will be awarded on style (5/8) "the code is clear to read" and efficiency (3/8) "your program is well constructed and run efficiently". For style, think about use of

CITS 1401 Computational Thinking with Python



proper comments, function docstrings, sensible variable names, your name and student ID at the top of the program, etc. (Please watch the lectures where this discussed).

Style Rubric:

0	Gibberish, impossible to understand.
1-2	Style is really poor or fair.
3-4	Style is good or very good, with small lapses.
5	Excellent style, really easy to read and follow.

Your program will be traversing text files of various sizes (possibly including large csv files), so you need to minimise the number of times your program looks at the same data items.

Efficiency rubric:

0	Code too complicated to judge efficiency or wrong problem tackled.
1	Very poor efficiency, additional loops, inappropriate use of <code>readline()</code> , etc.
2	Acceptable or good efficiency with some lapses.
3	Excellent efficiency, should have no problem on large files, etc.

Automated testing is being used so that all submitted programs are being tested the same way. Sometimes it happens that there is one mistake in the program that means that no tests are passed. If the marker can spot the cause and fix it readily, then they are allowed to do that and your - now fixed - program will score whatever it scores from the tests, minus 4 marks per intervention, because other students will not have had the benefit of marker intervention. Still, that's way better than getting zero. On the other hand, if the bug is hard to fix, the marker needs to move on to other submissions.

Extract from Australian Computing Society Accreditation manual 2019:

As per Seoul Accord section D, a complex computing problem will normally have some or all of the following criteria:

- involves wide-ranging or conflicting technical, computing, and other issues;
- has no obvious solution, and requires conceptual thinking and innovative analysis to formulate suitable abstract models;
- a solution requires the use of in-depth computing or domain knowledge and an analytical approach that is based on well-founded principles;
- involves infrequently encountered issues;
- is outside problems encompassed by standards and standard practice for professional computing;
- involves diverse groups of stakeholders with widely varying needs;
- has significant consequences in a range of contexts;
- is a high-level problem possibly including many component parts or sub-problems;
- identification of a requirement or the cause of a problem is ill defined or unknown.

Necessary formulas

i. **Correlation coefficient:**

Mathematical formula to calculate correlation is as follows:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

where x_i and y_i are the values of age and income of a specific user base respectively. \bar{x} is the average age and \bar{y} is the average income.

ii. **Standard deviation:**

$$s = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1}}$$

where $x_1, x_2, x_3 \dots x_n$ are observed values in the sample data. \bar{x} is the average value of observations and N is the number of observations.

data points