

Repeat Project 2018

Programming and Scripting

Due: last commit on or before August 22nd

This document contains the instructions for Repeat Project 2018 for Programming and Scripting. Note that you are not expected to know how to do the whole project from the beginning. Rather, it is expected that you will research ways to tackle the project and then formulate your own submission based on your investigations. Please be advised that all students are bound by the Quality Assurance Framework [2] at GMIT which includes the Code of Student Conduct and the Policy on Plagiarism.

Problem statement

The following project concerns the well-known BUPA Liver Disorders data set [3]. The project entails you researching the data set, and then writing documentation and code in the Python programming language [1] based on that research.

An online search for information on the data set will convince you that many people have investigated and written about it previously, and many of those are not experienced programmers. You are expected to be able to break this project into several smaller tasks that are easier to solve, and to plug these together after they have been completed. You might do that for this project as follows:

1. Research background information about the data set and write a summary about it.
2. Keep a list of references you used in completing the project.
3. Download the data set and write some Python code to investigate it.
4. Summarise the data set by, for example, calculating the maximum, minimum and mean of each column of the data set. A Python script will quickly do this for you.
5. Write a summary of your investigations.
6. Include supporting tables and graphics as you deem necessary.

It might help to suppose that your manager at work has asked you to investigate this data set, with a view to explaining it to your work colleagues. Imagine that you are to give a brief presentation on the data set to them in a few weeks' time, where you explain to them what investigating a data set entails and how Python can be used to do it. You have been asked not to create a deck of presentation slides, but rather to present your write-up and code to them.

Minimum Viable Project

The minimum standard for this project is a GitHub repository containing a README and a Python script. The README should contain a summary of the data set and your investigations into it. It should also clearly document how to run the Python code you used to investigate the data set, and what that code does. Furthermore, it should list all references used in completing the project.

A better project will be well organised and contain detailed explanations. The analysis will be well conceived, and examples of interesting analyses that others have pursued based on the data set will be discussed.

Note that the point of this project is to use Python. You may use any Python libraries that you wish, whether they have been discussed in class or not. You should not be thinking of using spreadsheet software like Excel to do your calculations.

Submissions

GitHub must be used to manage your project submission. Your GitHub repository will form the main submission of the project. You must submit the URL of your GitHub repository by email to ian.mcloughlin@gmt.ie before the deadline. You can do this at any time, as the last commit before the deadline will be used as your submission for this project.

Any submission that does not have a full and incremental git history with informative commit messages over the course of the project timeline will be accorded a proportionate mark. It is expected that your repository will have lots of commits, with each commit relating to a reasonably small unit of work. You may be asked at any time by the lecturer to explain the contents of your git repository. While it is encouraged that students will engage in peer learning, any unreferenced documentation and software that is contained in your submission must have been written by you. You can show this by having an incremental commit history and by being able to explain your code.

Marking scheme

The following marking scheme will be used to mark the project out of 100%. Students should note, however, that in certain circumstances the examiner's overall impression of the project may influence marks in each individual component.

25%	Research	Investigation of problem and possible solutions.
25%	Development	Clear analysis and well-written supporting code.
25%	Consistency	Good planning and pragmatic attitude to work.
25%	Documentation	Detailed descriptions and explanations.

Advice for students

- Your git commit history should be extensive. A reasonable unit of work for a single commit is a small function, or a handful of comments, or a small change that fixes a bug. If you are well organised you will find it easier to determine the size of a reasonable commit, and it will show in your git history.
- Using information, code and data from outside sources is sometimes acceptable — so long as it is licensed to permit this, you clearly reference the source, and the overall project is substantially your own work. Using a source that does not meet these three conditions could jeopardise your mark.
- You must be able to explain your project during and after its completion. Bear this in mind when you are writing your README. If you had trouble understanding something in the first place, you will likely have trouble explaining it a couple of weeks later. Write a short explanation of it in your README, so that you can jog your memory later.
- Everyone is susceptible to procrastination and disorganisation. You are expected to be aware of this and take reasonable measures to avoid them. The best way to do this is to draw up an initial straight-forward project plan and keep it updated. You can show the examiner that you have done this in several ways. The easiest is to summarise the project plan in your README. Another way is to use a to-do list like GitHub Issues.
- Students have problems with projects from time to time. Some of these are unavoidable, such as external factors relating to family issues or illness. In such cases allowances can sometimes be made. Other problems are preventable, such as missing the submission deadline because you are having internet connectivity issues five minutes before it. Students should be able to show that up until an issue arose they had completed a reasonable and proportionate amount of work and took reasonable steps to avoid preventable issues.
- Go easy on yourself - this is one project in one module. It will not define you or your life. A higher overall course mark should not be determined by a single project, but rather your performance in all your work in all your modules. Here, you are just trying to demonstrate to yourself, to the examiners, and to prospective future

employers, that you can take a reasonably straight-forward problem and solve it within a few weeks.

References

- [1] Python Software Foundation. Welcome to python.org.
<https://www.python.org/>.
- [2] GMIT. Quality assurance framework.
<https://www.gmit.ie/general/quality-assurance-framework>.
- [3] UC Irvine Machine Learning Repository. Liver disorders data set.
<https://archive.ics.uci.edu/ml/datasets/liver+disorders>.