



Practical Data Science with Python
COSC 2670/2738
Assignment 3

	Assessment Type	Individual
	Due Date	23:59 on the 15th of June, 2025
	Marks	30

Please read this carefully before attempting

This is an *individual* assignment. You may not collude with any other people, or plagiarise their work. You are expected to present the results of your own thinking and writing. Never copy other student's work (even if they "explain it to you first") and never give your written work to others. Keep any conversation high-level and never show your solution to others. Never copy from the Web or any other resource. Remember you are meant to generate the solution to the questions by yourself. Suspected collusion or plagiarism will be dealt with according to RMIT policy.

In the submission (your PDF file) you will be required to certify that the submitted solution *represents your own work only* by agreeing to the following statement:

I certify that this is all my own original work. If I took any parts from elsewhere, then they were non-essential parts of the assignment, and they are clearly attributed in our submission. I will show I agree to this honor code by typing "Yes":

Introduction

In this assignment, you are given a specific data science problem and several potential solutions. You are required to implement these solutions, then compare them and present critical analysis about how these techniques tackle the given data science problem.

The "Practical Data Science" Canvas contains further announcements and a discussion board for this assignment. Please be sure to check these on a regular basis – it is your responsibility to stay informed with regards to any announcements or changes. Login through <https://rmit.instructure.com/>.

Where to Develop Your Code

You are encouraged to develop and test your code in two environments: **Jupyter Notebook on Lab PCs** and **Anaconda 3** that was suggested in the course canvas announcement.

Jupyter Notebook on Lab PCs

On Lab Computer, you can find Jupyter Notebook via:

Start → All Programs → Anaconda3 (64-bit) → Jupyter Notebook

Then,

- Select New → Python 3
- The new created '*.ipynb' is created at the following location:
 - C:\Users\sXXXXXXX
 - where sXXXXXXX should be replaced with a string consisting of the letter "s" followed by your student number.

Academic integrity and plagiarism (standard warning)

Academic integrity is about honest presentation of your academic work. It means acknowledging the work of others while developing your own insights, knowledge and ideas. You should take extreme care that you have:

- Acknowledged words, data, diagrams, models, frameworks and/or ideas of others you have quoted (i.e. directly copied), summarised, paraphrased, discussed or mentioned in your assessment through the appropriate referencing methods
- Provided a reference list of the publication details so your reader can locate the source if necessary. This includes material taken from Internet sites. If you do not acknowledge the sources of your material, you may be accused of plagiarism because you have passed off the work and ideas of another person without appropriate referencing, as if they were your own.

RMIT University treats plagiarism as a very serious offence constituting misconduct. Plagiarism covers a variety of inappropriate behaviours, including:

- Failure to properly document a source
- Copyright material from the internet or databases
- Collusion between students

For further information on our policies and procedures, please refer to the following:

<https://www.rmit.edu.au/students/student-essentials/rights-and-responsibilities/academic-integrity>.

General Requirements

This section contains information about the general requirements that your assignment must meet. *Please read all requirements carefully before you start.*

- You *must* include a plain text file called “readme.txt” with your submission. This file should include your name and student ID, and instructions for how to execute your submitted script files. This is important as *automation* is part of the 6th step of data science process, and will be assessed strictly.
- Please ensure that your submission follows the file naming rules specified in the tasks below. File names are case sensitive, i.e. if it is specified that the file name is **gryphon**, then that is exactly the file name you should submit; **Gryphon**, **GRYPHON**, **griffin**, and anything else but **gryphon** will be rejected.

Overview

First, let's define $r_{a,i}$ as the rating from user u_a on item t_i , where u_a denotes the user whose index is a and t_i denotes the item whose index is i .

The core problem of recommender systems is to estimate how much the active user u_a would like a given target item t_i , typically expressed as a predicted rating $\hat{r}_{a,i}$.

Over the years, researchers and practitioners have proposed various recommendation algorithms. In this assignment, we focus on five selected algorithms, as outlined below:

- **Method 1:** Predicting $r_{a,i}$ as the average ratings of the corresponding users. Namely,

$$\hat{r}_{a,i} = \bar{u}_a, \quad (1)$$

where \bar{u}_a denotes the average rating of user u_a .

- **Method 2:** Predicting $r_{a,i}$ as the average ratings of the corresponding items. Namely,

$$\hat{r}_{a,i} = \bar{t}_i, \quad (2)$$

where \bar{t}_i denotes the average rating of item t_i .

- **Method 3:** Predicting $r_{a,i}$ by using the User KNN-based Collaborative Filtering. Namely, $r_{a,i}$ is predicted by using the User KNN-based Collaborative Filtering. However, if all of the k nearest neighbours have no ratings on item t_i , $r_{a,i}$ will be predicted as the average rating of user u_a .
- **Method 4:** Predicting $r_{a,i}$ by using the Item KNN-based Collaborative Filtering. Namely, $r_{a,i}$ is predicted by using the Item KNN-based Collaborative Filtering. However, if all of the k nearest neighbours have no ratings from u_a , $r_{a,i}$ will be predicted as the average rating of item t_i .
- **Method 5:** Predicting $r_{a,i}$ by using both the user KNN-based Collaborative Filtering and the Item KNN-based Collaborative Filtering. Namely, $r_{a,i}$ is predicted as

$$\hat{r}_{a,i} = \lambda \hat{r}_{a,i}^u + (1 - \lambda) \hat{r}_{a,i}^t, \quad (3)$$

where $\hat{r}_{a,i}^u$ denotes the prediction of $r_{a,i}$ by using User KNN-based Collaborative Filtering, $\hat{r}_{a,i}^t$ denotes the prediction of $r_{a,i}$ by using Item KNN-based Collaborative Filtering, and λ is a pre-defined parameter for how to integrate the predictions from $\hat{r}_{a,i}^u$ and $\hat{r}_{a,i}^t$, and the value of λ is between 0 and 1.

Tasks

Task 1: Implementation

In this task, you are required to implement the above five methods (solutions) within the scenario configured in the provided *assignment3_framework.ipynb*. Please note that: for Method 3, 4 and 5, you are required to utilize the knowledge covered in this course to optimize their performance.

This Python framework file, named *assignment3_framework.ipynb*, is designed to help you get started and will also automate the correctness marking process. The framework includes both the training data and the test data.

Please read the comments in the provided assignment framework carefully, and insert your own code only in the designated code cells as indicated.

Please do not change anything else in the remaining cells of the framework, as doing so may result in errors during the automatic marking process, which might cause the submission invalid.

Please provide detailed comments to explain your implementation. To what level of details should you provide in your solution? Please take the comments in the *ipynb* files in Week 10 (*knn_based_cf_updated.zip*) as examples for the level of detailed comments you are expected to put for your solution. You might find the following information useful: https://www.w3schools.com/python/python_comments.asp

Note, you are required to implement your own implementation, and please do not use any other libraries that are related to Recommender Systems or Collaborative Filtering. If you use any of these libraries, your implementation part will be invalid.

Task 2: Presentation

- The presentation should
 - Explain how Method 3, 4 and 5 are optimized to achieve their best performance.
 - Compare and analyse the performance of all these 5 methods.
 - Explain why some methods are better than others.
- The presentation should be no more than 10 minutes.
- Your presentation slides should be:
 - Microsoft PowerPoint slides (with audio inserted for each slide by using: **Insert** – > **Audio** – > **Record Audio**).
 - or you can create your own presentation slides (e.g. PDF version) and please submit your own recording (in the format of mp4 or avi) of your presentation as well.

What to Submit, When, and How

The assignment is due at

23:59 on the 15th of June, 2025.

Assignments submitted after this time will be subject to standard late submission penalties.

The following files should be submitted:

- Notebook file containing your python implementation, 'Assignment3_framework.ipynb'.

For the notebook file, follow these steps before submission:

1. Main menu → Kernel → Restart & Run All
2. Wait till you see the output displayed properly. You should see all the data printed and graphs displayed.

- One of the following:
 - Your **Slides.pdf** file and your presentation recording in the required format.
Or,
 - Your Microsoft PowerPoint slides (with audio inserted for each slide).
- The “readme.txt”: includes your name and student ID, and instructions for how to execute your submitted script files.
- Please note: there is no need to submit the data sets, as you are not allowed to change them.

They must be submitted as ONE single zip file, named as your student number (for example, 1234567.zip if your student ID is s1234567). The zip file must be submitted in Canvas:

Assignments/Assignment 3.

Please do NOT submit other unnecessary files.