

School of Computing and Information Systems
The University of Melbourne
COMP30027, Machine Learning, 2023

Project 2: Book Rating Prediction

Task:	Build a classifier to predict the rating of books
Due:	Group Registration: Friday 5 May, 5pm Stage I: Friday 19 May, 5pm Stage II: Friday 26 May, 5pm
Submission:	Stage I: Report (PDF) and code to Canvas; test outputs to Kaggle in-class competition Stage II: Peer reviews and reflection to Canvas
Marks:	The Project will be marked out of 20, and will contribute 20% of your total mark.
Groups:	Groups of 1 or 2, with commensurate expectations for each (see Sections 2 and 5).

1 Overview

The goal of this Project is to build and critically analyse supervised Machine Learning methods to predict the ratings of books based on their titles, authors, descriptions and other features. There are three levels of rating, 3, 4 and 5, for each book.

This assignment aims to reinforce the largely theoretical lecture concepts surrounding data representation, classifier construction, evaluation and error analysis, by applying them to an open-ended problem. You will also have an opportunity to practice your general problem-solving skills, written communication skills, and critical thinking skills.

2 Deliverables

This project has two stages. The deliverables of each stage are listed as follows. More details about deliverables are given in the Submission (Section 5).

Stage I:

1. **Report:** an **anonymous** written report, of 1,300-1,800 words (for a group of one person) or 2,000-2,500 words (for a group of two people).
2. **Output:** the output of your classifiers, comprising the label predictions for test instances, submitted to the Kaggle¹ in-class competition described below.
3. **Code:** one or more programs, written in Python, which implement machine learning models to make predictions and evaluate the results.

Stage II:

1. **Peer review:** reviews of two reports written by other students, 200-300 words each (for a group of one person) or 300-400 words each (for a group of two people).
2. **Reflection:** a written reflection piece of 400 words. This deliverable is individual work.

¹<https://www.kaggle.com/>

3 Data

The information of book is collected from Goodreads², which is a platform that allows users to search its database of books, rate books and write reviews. The data files for this project are available via Canvas, and are described in a corresponding README.

In our dataset, each book contains:

• **Book features:** name, authors, publish year, publish month, publish day, publisher, language, page numbers, and description.

• **Text features:** produced by various text encoding methods for name, authors, and description. Each text feature is provided as a single file with rows corresponding to the file of book features.

• **Class label:** the rating of a book `rating_label` (3 possible levels, 3, 4 or 5)

ordinal
not categorical

You will be provided with a training set and a test set. The training set contains the book features, text features, and the `rating_label`, which is the “class label” of our task. The test set only contains the book and text features without labels.

The files provided are:

- *book_rating_train.csv*: the book features and class label of training instances.
- *book_rating_test.csv*: the book features of test instances.
- *book_text_features_*.zip*: the preprocessed text features for training and test sets, 1 zipped file for each text encoding method. Details about using these text features are provided in README.

4 Task

You are expected to develop Machine Learning models to **predict the rating of a book based on its features** (e.g. name, authors, description, publish year etc.). You will **explore effective features, implement and compare different machine learning models and conduct error analysis** for this task.

Various machine learning techniques have been (or will be) discussed in this subject (**OR, Naive Bayes, Decision Trees, kNN, SVM, neural network, etc.**); many more exist. You may use any machine learning method you consider suitable for this problem. *You are strongly encouraged to make use of machine learning software and/or existing libraries (such as `sklearn`) in your attempts at this project.*

In addition to different learning algorithms, there are many different ways to **encode text** for these algorithms. The files in *book_text_features_*.zip* are some possible representations of the name, authors and description of books we have provided. For example, one of the encoding method is **CountVectorizer in sklearn**, which converts text documents into “Bag of Words” – the documents are described by word occurrences while ignoring the relative position information of the words. You can use these representations to develop your classifiers, but please also feel free to extract your own features from the raw book features according to your needs. Just keep in mind that any data representation you use for the text in the training set will need to be able to generalise to the test set.

other
encoder

You are expected to complete the following two phases for this task:

- **Training-evaluation phase:** the **holdout** or **cross-validation** approaches can be applied on the training data provided.
- **Test phase:** the trained classifiers will be evaluated on the unlabelled test data. The predicted labels of test cases should be submitted as part of the Stage I deliverable.

²<https://www.goodreads.com/>

5 Submission

The report, code, peer reviews and reflections should be submitted via Canvas; the predictions on test data should be submitted to Kaggle.

5.1 Individual vs. Team Participation

You have the option of participating individually, or in a group of two. In the case that you opt to participate **individually**, you will be required to **implement at least 2 and up to 4** distinct Machine Learning models. **Groups of two** will be required to **implement at least 4 and up to 5** distinct Machine Learning models, of which **one is to be an ensemble model – stacking based on the other models**. The report length requirement also differs, as detailed below:

Group size	Distinct models required	Report length
1	2–4	1,300–1,800 words
2	4–5	<u>2,000–2,500</u> words

Group Registration

If you wish to form a group of 2, **only one** of the members needs to register by **Friday 5 May 5:00pm**, via the form “**Project 2 Group Registration**” on Canvas. For a group of 2, **only one** of the members needs to submit deliverables.

Note that once you have signed up for a given group, you will not be allowed to change groups. If you do not register before the deadline above, we will assume that you will be completing the assignment as an individual, even if you were in a two-person group for Assignment 1.

5.2 Stage I: Report

Your report is expected to demonstrate the knowledge that you have gained and the critical analysis you have conducted in a manner that is accessible to a reasonably informed reader.

The report should be 1,300-1,800 words (individual) or 2,000-2,500 words (groups of two people) in length **excluding reference list, figure captions and tables**. The report should include the following sections:

1. **Introduction**: a basic description of the task and a short summary of your report.
2. **Methodology**: what you have done, including any learners that you have used, and **features that you have engineered**. *This should be at a conceptual level; a detailed description of the code is not appropriate for the report. The description should be similar to what you would see in a machine learning conference paper.*
3. **Results**: performance of your classifiers, in terms of **evaluation metric(s)** and, ideally include **figures** and **tables**.
4. **Discussion and Critical Analysis**: this section should include a more *detailed discussion* which contextualises the **behaviour of the method(s)**, in terms of the **theoretical properties** we have identified in the lectures and **error analysis of the method(s)**. This is the most important section of the report.
5. **Conclusion**: demonstrates your identified knowledge about the problem.
6. Reference: reference of related work.

Note that we are more interested in seeing evidence that you have thought about the task and investigated the reasons for the relative performance of different methods, rather than in the raw accuracy/scores of different methods. This is not to say that you should ignore the relative performance of different runs over the data, but

rather than you should think beyond simple numbers to the reasons that underlie them, and connect these to the theory that we have discussed in this subject.

We provide L^AT_EX and Word style files that we would prefer that you use in writing the report. Reports must be submitted in the form of **a single PDF file**. If a report is submitted in any format other than PDF, we reserve the right to return the report with a mark of 0.

To facilitate anonymous peer-review, your name and student ID should not appear anywhere in the report, including the metadata (filename, etc.).

5.3 Stage I: Predictions of test data

To give you the possibility of evaluating your models on the test set, we will be setting up a Kaggle in-class competition for this project. You can submit results on the test set there, and get immediate feedback on your model's performance. There is a Leaderboard, that will allow you to see how well you are doing as compared to other classmates participating online. The Kaggle in-class competition URL and instructions will be announced on Canvas shortly.

You will receive marks for submitting at least one set of predictions for the unlabelled test set into the competition; and get basically **reasonable accuracy (higher than 60%)**. The focus of this assignment is on the quality of your critical analysis and your report, rather than the performance of your Machine Learning models.

5.4 Stage II: Reviews

Stage II submissions will be open as soon as the reports are available (within 24 hours after the Stage I submission deadline). During the reviewing process, you will read two submissions by other students. This is to help you contemplate some other ways of approaching the project, and to ensure that students get some extra feedback. For each report, you should aim to write 200-400 words total (200-300 words if you work alone, or 300-400 words if you work in a group of two people), responding to three questions:

- Briefly summarise what the author has done
- Indicate what you think that the author has done well, and why
- Indicate what you think could have been improved, and why

Please be courteous and professional in the reviewing process. A brief guideline for reviewers published by IEEE can be found https://www.ieee.org/content/dam/ieee-org/ieee/web/org/members/students/reviewer_guidelines_final.pdf.

5.5 Stage II: Reflections

A comprehensive written reflection piece summarising your critical reflection on the following topics **within 400 words**.

1. the **process of completing this project**
2. things that you **are satisfied with** and those **can be improved** in your Stage I deliverables, e.g. modelling, evaluation, analysis and discussion.
3. if you have worked in a group of two people, what are the **individual contributions**?

The reflection report is **individual and not anonymous**. Everyone must submit their own reflection on Canvas.

6 Assessment Criteria

The Project will be **marked out of 20**, and is worth 20% of your overall mark for the subject. The mark breakdown will be:

Stage I	Report	15 marks
	Performance of classifier	2 mark
Stage II	Reviews	2 marks
	Reflection	1 mark
TOTAL		20 marks

The report will be marked according to the rubric, which is published on the Canvas. You have to submit your code that supports the results presented in your report. If you do not submit an executable code that supports your findings, you will receive a mark of 0 for the “Report”.

The performance of classifier is for **submitting at least one set of model predictions to the Kaggle competition** (1 mark); and getting basically reasonable accuracy (**higher than 60%**) (1 mark).

Since all the data exist on the World Wide Web, it is possible to “cheat” and identify some of the class labels from the test data using non-Machine Learning methods. If there is any evidence of this, the performance of classifier will be ignored, and you will instead receive a mark of 0 for the “Performance of classifier”.

7 Using Kaggle

The Kaggle in-class competition URL will be announced on Canvas shortly. To participate the competition:

- Each student should create a Kaggle account (unless they have one already) using your Student-ID
- You may make up to 8 submissions per day. An example submission file can be found on the Kaggle site.
- Submissions will be evaluated by Kaggle for accuracy, **against just 50% of the test data**, forming the public leaderboard.
- Prior to competition close, you may select a final submission out of the ones submitted previously – by default the submission with highest public leaderboard score is selected by Kaggle.
- After competition close, public test scores will be replaced with the private leaderboard test scores (100% test data).

8 Assignment Policies

8.1 Terms of Use

Please note that the dataset is a sample of actual data posted to the World Wide Web. As such, it may contain information that is in poor taste, or that could be considered offensive. We would ask you, as much as possible, to look beyond this to the task at hand. For example, it is generally not necessary to read individual records.

The opinions expressed within the data are those of the anonymised authors, and in no way express the official views of the University of Melbourne or any of its employees; using the data in an educative capacity does not constitute endorsement of the content contained therein.

If you object to these terms, please contact Ling Luo (ling.luo@unimelb.edu.au) as soon as possible.

8.2 Changes/Updates to the Assignment Specifications

We will use Canvas to advertise any (hopefully small-scale) changes or clarifications in the assignment specifications. Any addenda made to the assignment specifications via Canvas will supersede information contained in this version of the specifications.

8.3 Late Submissions

Late submissions will bring disruption to the reviewing process. You are strongly encouraged to submit by the date and time specified above. If circumstances do not permit this, then the marks will be adjusted as follows:

- Each day (or part of the day) that the report is submitted after the specified due date and time for Stage I, 10% will be deducted from the marks available, up until 7 days (1 week) has passed, after which regular submissions will no longer be accepted. A late report submission will mean that your report might not participate in the reviewing process, and so you will probably receive less feedback.
- Any late submission of the reviews will incur a 50% penalty (i.e. 1 of the 2 marks), and will not be accepted more than 7 days (1 week) after the Stage II deadline.
- Any late submission of the reflection will incur a 50% penalty (i.e. 0.5 of the 1 mark), and will not be accepted more than 7 days (1 week) after the Stage II deadline.

8.4 Academic Honesty

While it is acceptable to discuss the assignment with others in general terms, excessive collaboration with students outside of your group is considered collusion. Your submissions will be examined for originality and will invoke the University's Academic Misconduct policy (<https://academicintegrity.unimelb.edu.au/>) where either inappropriate levels of collaboration or plagiarism are deemed to have taken place.