# Book Recommendation System Project

# Letter of Transmittal and Project Proposal

## Letter of Transmittal

September 7, 2024

Mr. John Smith

XYZ Book Company

505 Central Ave

Pacific Grove, California

Dear Mr. Smith,

Given the increasing growth of online book sales, more and more organizations are developing book recommendation systems. Such systems work with data to recommend book titles that align with their customer's interests. These recommendations drive sales and enhance the customer experience by aiding the discovery of new titles that interest them. Having identified what a book recommendation system is, we can move on to discuss the current problem XYZ Book Company faces and why the development of a recommendation system would be beneficial to the company.

XYZ Book Company's platform currently lacks any ability to suggest books when customers search for a title. This prevents the customer from discovering additional titles that they might be interested in, potentially resulting in reduced customer satisfaction when compared to other platforms that do have the ability to recommend titles. Reduced customer satisfaction will negatively impact our business; thus, it is important that we respond to this problem quickly and find a solution.

Our current proposed solution is a user-friendly application that lets customers enter a book title as text. This text is analyzed by the application, and titles that are similar to the input title will be recommended. This implementation will allow customers to discover new books that align with their interests, enhancing the user experience on our platform. Providing suggestions in this way will increase the engagement of our customers, encourage exploration of recommended titles, and ultimately lead to higher sales.

The costs associated with the application are an initial $32,160 and $720 yearly for maintenance. The time it would take to build the application should be no longer than eighteen days. Gathering requirements will take two days. Creating a design plan based on those requirements will take five days. Processing the data used in the application will take two days. Creating and training

the machine learning model will take two days. Developing the user interface will take one day. Testing the application will take three days. Finally, deploying the application will take three days. The data used for the application will be sourced from Kaggle.com and includes customer ratings data for individual book titles, which will allow our application to make suggestions based on customer tastes.

My personal experience includes creating multiple applications that implement machine learning, as well as many other projects that work with data. I also have relevant education with a bachelor's degree in computer science from Western Governors University. With this expertise, I trust in my ability to successfully create an application that will solve the problem XYZ Book Company faces and benefit the company.

Thank you for your time. I look forward to hearing your response.

Sincerely,

*Samuel Buck*

Samuel Buck, Software Developer

# Project Proposal

## Project Summary

XYZ Book Company does not currently have a way to recommend books to its customers based on their interests. This limits the customer's ability to search for books on our platform, which hinders customer engagement and satisfaction. To rectify this incapacity, we intend to build a book recommender system using a machine-learning model to suggest books similar to what customers search for on our platform. The model will work with a set of data containing user rating information on individual book titles, enabling the model to predict titles similar to a given title. This capability will benefit XYZ Book Company by increasing sales, as customers are more likely to purchase titles that align with their interests.

The methodology we will use to implement this project will be the waterfall method. This methodology establishes a sequential approach and is suitable for a project that has concrete requirements.

This project will have six phases that are listed below:

1. Requirements Gathering
   We will document XYZ Book Company's requirements and define the scope of the project. During this phase, the necessary tools and data required to develop the recommender system will be identified.

2. System Design
   We will create a design plan for the application's architecture, including the final dataset, the machine learning algorithm, and the user interface. This plan will benefit the development of the system by establishing guidelines for the design.

3. Implementation
   Using the design plan created in the previous phase, we will process data to create a prepared dataset, create and train the machine learning model on the dataset, and develop the user interface.

4. Testing
   After implementing the design plan, we will test the machine learning model's performance, as well as the user interface. When these results are deemed satisfactory, we will move on to deploy the application in the next phase.

5. Deployment
   We will deploy the system to an environment that can be used by XYZ Book Company. Setting up the necessary infrastructure, packaging the application, and preparing a user guide will all be completed during this phase.

6. Maintenance

The final maintenance phase will constitute the monitoring of the system and providing support to address any issues that present themselves. User feedback may also be gathered in this phase to support the improvement of the application.

# Executive Summary

## Project Summary

XYZ Book Company currently lacks any way to recommend books to its customers on its platform. This limits the engagement of customers, which may result in missed opportunities to increase sales and achieve a higher level of customer satisfaction. XYZ Book Company is an online book retailer that aims to improve the customer experience by introducing a book recommendation system to its platform. This recommendation system will use a nearest neighbors machine learning model to predict similar titles. The model will be trained on a dataset containing user rating data for individual books, allowing for a reasonably accurate method of suggesting titles based on an input title.

We will provide two deliverables: the application itself and a user guide to aid in setting up the application. The application is a book recommendation system that suggests book titles similar to an input title. It will utilize a collaborative filtering approach with the use of a k-nearest neighbors algorithm to analyze patterns in user ratings and identify titles that may interest these users. When a user inputs a title, the system will locate the book within the dataset and then find the nearest neighbors using a machine learning model. These neighbors will then be presented to the user as recommendations. The user will interact with the user interface to enter a book title, submit the title, and receive recommendations based on the input title.

The user guide will provide a step-by-step guide to setting up the application on a local system. The user guide will encompass the installation of prerequisite items, instructions for setting up an environment, and information on how to launch and use the application. The guide will also provide information for troubleshooting errors that may occur during setup or use of the application.

The recommendation system will enhance XYZ Book Company's customer experience by providing suggestions based on user preferences. These recommendations will help customers discover new titles that will likely interest them. This will benefit the company by increasing sales, as users who are suggested books that are pertinent to their tastes will be more likely to purchase them. Ultimately, by establishing this recommender system in XYZ Book Company's platform, we will benefit from an improved shopping experience and increased revenue.

## Data Summary

The data used in the application will be sourced from Kaggle.com and will be in the form of two datasets. The first dataset will contain user rating data, such as the user's ID, the rated book's ISBN, and the rating the user gave. The second dataset will contain book data, such as the book's ISBN, title, author, publication year, and publisher. These two datasets will provide the necessary foundation to build a recommendation system.

The data will be downloaded directly from Kaggle and imported into the application for analysis. During the design phase, the data will be explored to detect patterns and opportunities for feature engineering. The data will also be cleaned in this phase, constituting the handling of any missing values and outliers. The development phase will involve transforming data into a format suitable for a machine learning algorithm, specifically the k-nearest neighbor algorithm from the Scikit-Learn Python library.

The data will meet the needs of the project, as it includes a wide variety of user ratings and book data, thus making the data suitable for use in our recommender system. Data anomalies, such as rows with missing data or outliers, will be addressed during the initial design phase. Rows that contain missing or incorrect data will be corrected if possible or removed from the dataset to ensure that the model is not negatively impacted. Additionally, entries such as books with too few ratings or users with too few ratings will be excluded to ensure they do not impact the machine learning model.

As it is important for the recommender system to be free of any ethical or legal concerns, precautions have been taken to ensure that the dataset used for the application will meet our expectations for these considerations. Because the user data will be anonymized and the book data is publicly available information, there are minimal ethical concerns. As for legal concerns, the dataset is licensed under CC0 1.0, allowing free use even for commercial purposes.

## Implementation

This project will be implemented using the waterfall methodology. Using this methodology, we will have a sequential approach where each project phase must be completed before moving on to the next. The waterfall methodology is suitable for projects with concrete requirements and will allow us to plan and execute our project in a linear fashion.

The phases for this project using the waterfall methodology are as follows:

1. Requirements Gathering
   In the requirements gathering phase, we will collect and document all requirements from XYZ Book Company. We will define the scope of the project, as well as identify the tools and data necessary to create the recommender system.

2. System Design
   In the system design phase, we will plan out the architecture of the application, involving the final dataset, the applied machine learning algorithm, and the user interface. This design plan will aid in the development of the recommender system, allowing us to focus on implementing designs that have been previously established.

3. Implementation
   The system will be developed using the design plan established in the previous phase. This phase will constitute the processing of data, creation of the nearest neighbors machine learning model, and development of a user interface. The machine learning

model will also be trained on the prepared dataset in this phase.

4. Testing
We will measure the model's performance in this phase using precision and recall. Testing of the user interface and other application features will be performed during this phase. When we are satisfied with the results of the testing phase, we will move on to the next phase.

5. Deployment
When the testing phase is completed, we will deploy the system to an environment that is usable by XYZ Book Company. This will involve setting up the necessary infrastructure, packaging the application, and preparing the user guide.

6. Maintenance
In the maintenance phase, we will monitor the system and provide support to address any issues that may arise. We may also update the application to improve functionality based on user feedback.

# Timeline

We anticipate the timeline of this project to be between 9/9/24 and 10/2/24, taking a total of 18 business days to complete. A table containing a timeline for deliverables and milestones can be found below.

| Milestone | Duration (days) | Projected start date | Anticipated end date |
|---|---|---|---|
| Requirements Gathering | 2 | 9/9/24 | 9/10/24 |
| Design Plan | 5 | 9/11/24 | 9/17/24 |
| Process Data | 2 | 9/18/24 | 9/19/24 |
| Create and Train Machine Learning Model | 2 | 9/20/24 | 9/23/24 |
| User Interface | 1 | 9/24/24 | 9/24/24 |

| | | | |
|---|---|---|---|
| Test Machine Learning Model and User Interface | 3 | 9/25/24 | 9/27/24 |
| Application Deployment | 3 | 9/30/24 | 10/2/24 |

## Evaluation Plan

To ensure that the application adheres to our requirement guidelines, the design plan will undergo peer reviews to ensure that it meets XYZ Book Company's requirements before moving on to the implementation phase. We will conduct code reviews during the implementation phase to ensure that the code written follows industry-standard best practices. Each application component will be tested to ensure that it functions as intended. This testing will include the data preprocessing scripts, machine learning model, and user interface. During the testing phase, we will measure the machine learning model's precision and recall to ensure that our model meets at least 70% precision and 50% recall. We will also perform user acceptance testing on the system to ensure that the application features and user interface meet our requirements. During the deployment phase, we will verify that the system is correctly installed and operational.

After deployment, the model's recommendations will be validated by gathering user feedback. Real-world validation will allow for adjustments and fine-tuning of the model to improve its accuracy and relevance in the future. This user feedback will be gathered by surveying users after the purchase of a book and asking their opinions on the titles that were recommended to them.

# Resources and Costs

## Hardware and Software Costs:

- **Development Workstations** - $1,200 (per employee)
- **Google Cloud Storage** - $10/month
- **Server** - $1,000
- **Miniconda Python Environment** – Free
- **Data Analytics Libraries –** Free
- **Jupyter Lab** – Free
- **Kaggle Dataset** - Free

Total: $4600 + $10/month

## Estimated Labor and Costs:

- **Project Manager -** $60/hour | ($8,640 total)
- **Data Scientist** - $70/hour | ($10,080 total)
- **Software Developer** - $60/hour | ($8,640 total)

Total: $27,360 for the estimated duration of the project

## Estimated Environment Costs:

- **Initial Setup -** $200
- **Monitoring and Maintenance -** $50/month

Total: $200 + $50/month

Total Resources and Costs Estimate: $32,160 + $60/month

# Application

## Submitted Files List

Listed below are the application files that will be submitted:

\recommender_application_files

    \environment.yml          Used to create the Miniconda environment.

    \book_recommender.ipynb     This is the notebook file that contains the application code.

    \csv_files

        \books.csv     This is a dataset containing books and book information.

        \ratings.csv     This is a dataset containing user ratings data for specific titles.

# Post-implementation Report

## Solution Summary

XYZ Book Company faced many issues with its lack of a recommendation system. The inability to suggest books to customers was a limiting factor in sales and customer engagement and allowed business competitors to increase their customer base due to our customers migrating to their platforms. This problem required the creation of an application that could provide recommendations to users based on an input title. The application achieved this through the usage of a machine learning model utilizing the k-nearest neighbors algorithm to identify the similarity between books based on user rating data.

The application makes use of user-based collaborative filtering to examine the preferences of a large number of users in the dataset. This involved the creation of a matrix to map individual user ratings to specific books, which would then be analyzed by the machine learning model to identify relationships between books. The two images below show the pivot table created from the final dataset and its mappings on the matrix.

| User-ID Book-Title | 254 | 507 | 882 | 1424 | 1435 |
|---|---|---|---|---|---|
| 1984 | 9.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1st to Die: A Novel | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2010: Odyssey Two | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 204 Rosewood Lane | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 24 Hours | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

```
[[9., 0., 0., ..., 0., 0., 0.],
 [0., 0., 0., ..., 0., 0., 0.],
 [0., 0., 0., ..., 0., 0., 0.],
 ...,
 [0., 0., 0., ..., 0., 0., 0.],
 [0., 0., 0., ..., 0., 0., 0.],
 [0., 0., 0., ..., 0., 0., 0.]]
```

When the application is given an input title, the application finds books that have similar user rating patterns. This allows users to receive recommendations for books that they enjoy based on the data from users who have similar literary interests established from their rating data. That is to say, if the input title and another title have both been rated highly by many users, then the likelihood that those two titles are similar is increasingly likely.

The k-nearest neighbors algorithm carried out the primary function of the system and was essential in providing the capability to recommend book titles when given user input. The algorithm identifies titles that are nearest to the given title by computing the cosine similarity between all books in the matrix. The application uses the book title as an identifier to find the corresponding index in the pivot table, then finds the k books with the smallest distance to the input book.

The application provides a way for users to receive recommendations based on an input title and is a valuable tool for XYZ Book Company to maintain customer satisfaction and engagement. As a result of this application, we can expect the number of purchases to increase, as it will be easier for customers to identify titles that match their interests. This application also increases the competitiveness of XYZ Book Company, expanding its share of the online literature market.
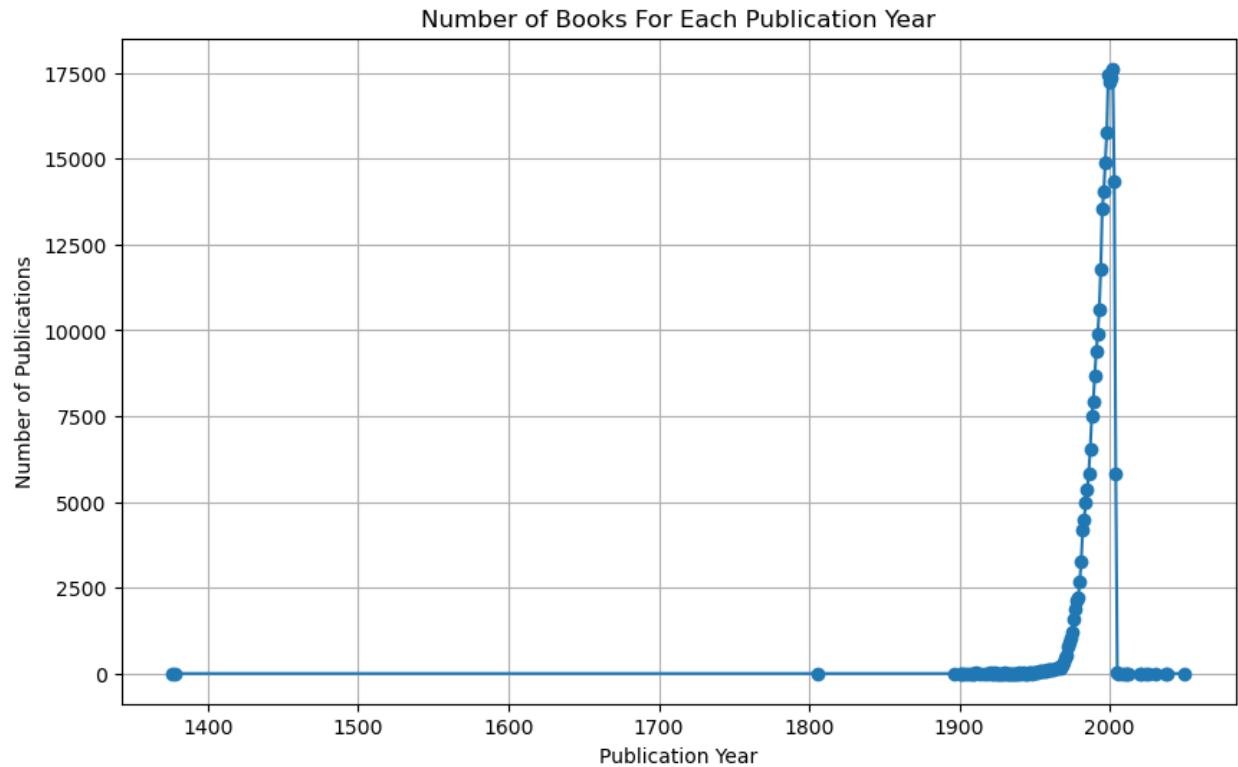
## Data Summary

The raw data for this project was sourced from Kaggle. The data obtained included a dataset containing individual book data and a dataset containing user ratings data. These datasets were in the comma-separated values format and stored in the application as data frames using the panda's library. The images below show the structure of the book data frame and the ratings data frame.

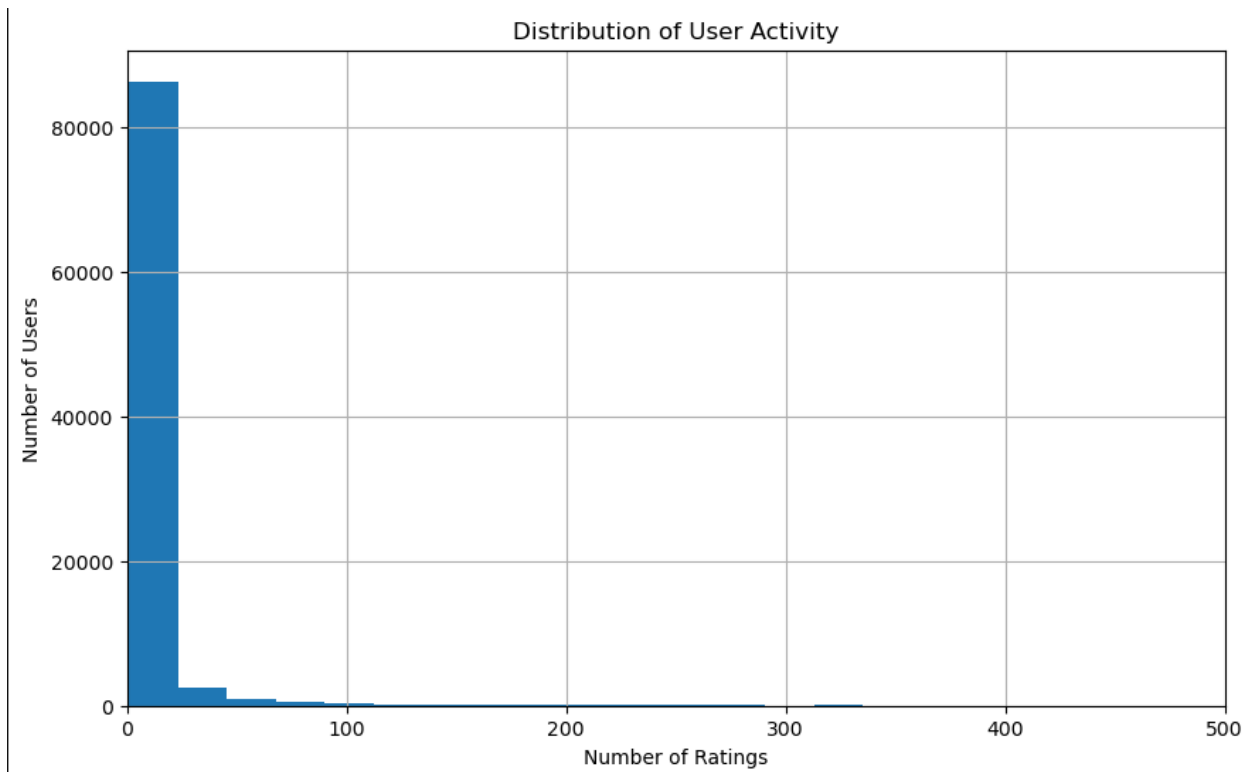| | ISBN | Book-Title | Book-Author | Year-Of-Publication | Publisher |
|---|---|---|---|---|---|
| 0 | 0195153448 | Classical Mythology | Mark P. O. Morford | 2002 | Oxford University Press |
| 1 | 0002005018 | Clara Callan | Richard Bruce Wright | 2001 | HarperFlamingo Canada |
| 2 | 0060973129 | Decision in Normandy | Carlo D'Este | 1991 | HarperPerennial |
| 3 | 0374157065 | Flu: The Story of the Great Influenza Pandemic... | Gina Bari Kolata | 1999 | Farrar Straus Giroux |
| 4 | 0393045218 | The Mummies of Urumchi | E. J. W. Barber | 1999 | W. W. Norton &amp; Company |

| | User-ID | ISBN | Book-Rating |
|---|---|---|---|
| 0 | 276725 | 034545104X | 0 |
| 1 | 276726 | 0155061224 | 5 |
| 2 | 276727 | 0446520802 | 0 |
| 3 | 276729 | 052165615X | 3 |
| 4 | 276729 | 0521795028 | 6 |

The dataset was examined during the system design phase to ensure it met the application requirements. It was found that the book data frame contained two missing 'Book-Author' values and two missing 'Publisher' values. Although these columns would not ultimately be used in the machine learning model, the values were filled with the string 'Unknown' to ensure no null values were in the dataset.

The number of books for each publication year was also examined in this phase. This information would help us understand the dataset and its makeup. It was found that the overwhelming majority of titles in the dataset were published before the year 2000. While the dataset lacks titles after the year 2000, this did not hinder the development of the application.

Number of Books For Each Publication Year

The final data frame was created by merging the user ratings data frame and the book data frame on the ISBN number. This combined the two data frames to create a data frame that contained user rating data as well as book data. The number of occurrences for each unique user ID in the data frame was counted to determine the distribution of user activity. The results showed that the majority of users rated very few books. These users, being very large in number and potentially affecting the machine learning model's accuracy, would need to be filtered out.

Distribution of User Activity

Users who rated fewer than 100 books were removed from the data frame to ensure that the machine-learning model performed well. Books that had fewer than 50 ratings were also removed from the data frame. Filtering this data ensured that the recommendation system would not be negatively affected by data of insufficient quantity and make more accurate predictions.

In the implementation phase, the final data frame was transformed into a pivot table. The pivot table rows would represent book titles, and the columns would represent users. The cells in the pivot table contained the user's rating for an individual book, and NA values were replaced with 0. This pivot table would then be converted into a matrix, which would be analyzed by the machine learning model. An image of the final pivot table can be found below.

| User-ID | 254 | 507 | 882 | 1424 | 1435 |
|---|---|---|---|---|---|
| **Book-Title** | | | | | |
| 1984 | 9.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1st to Die: A Novel | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2010: Odyssey Two | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 204 Rosewood Lane | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 24 Hours | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

The CSV files may be expanded to include new user ratings and book entries during the maintenance phase. The model can then be retrained on the data to ensure our recommender system stays relevant as new titles and user ratings are obtained. This supports the engagement of XYZ Book Company's customers with the application, as new titles may possibly be recommended to them.

## Machine Learning

The book recommendation system suggests books that are similar to a user's input by using an implementation of the k-nearest neighbors algorithm. Users enter a book title, and the system returns similar books based on the identified user ratings similarities that the algorithm has analyzed. By using k-nearest neighbors in our recommendation system, we are able to identify patterns within the user ratings data and discover relationships between different books. This gives the application the functionality to recommend books to a user and helps solve the many problems faced by XYZ Book Company's lack of a book recommendation system.

Our recommender system uses a user-based collaborative filtering approach and k-nearest neighbors to recommend books based on analyzing users with similar rating tendencies. This implementation uses the Scikit-Learn library's nearest neighbor function to compute the cosine similarity between all rating data points in the matrix, thus identifying similar books. The application also utilizes the Pandas library to import and process the data, as well as NumPy to create and work with matrices. Matplotlib was used to create visualizations from the data. The application was created using Jupyter Lab, enabling the construction of a single notebook that contains the proposed recommendation functionality. Miniconda was used to create the environment and enabled the straightforward installation of software dependencies.

The implementation plan included data preprocessing, model development, and model improvement. Missing values and sparse user data were removed from the dataset and ultimately transformed into a pivot table matrix. The matrix created from the pivot table would be used to train the machine learning model. The k-nearest neighbors model was created using the Scikit-Learn library and trained on the pivot table matrix. Each cell in the matrix represented a user's rating of a specific book, and the cosine distance between books was calculated by the model to identify the nearest neighbors for each book in the matrix. The minimum ratings count for users, and the book ratings count were adjusted to find the optimal number for use by the machine learning model. Precision and recall metrics were examined for different k values to test the model's performance. The resulting machine learning model was able to accurately predict book titles and could successfully be used in the recommendation system.

The model uses k-nearest neighbors due to its effectiveness in user-based collaborative filtering projects. The algorithm effectively handled the sparse data in the pivot table matrix and provided strong predictions based on user-book rating interactions. The training process involved the calculation of distances between points over the entire dataset, which could then be evaluated using precision and recall. This training method was chosen because the model is not expected to come across any unexpected data and is intended to make predictions based on all available user

ratings. Precision and recall metrics were used to evaluate the model's performance, as they provide valuable information about the quality of recommendations by measuring the overall relevancy of recommendations and the number of relevant books that were successfully recommended. These decisions ensured that the model effectively made relevant recommendations to users.
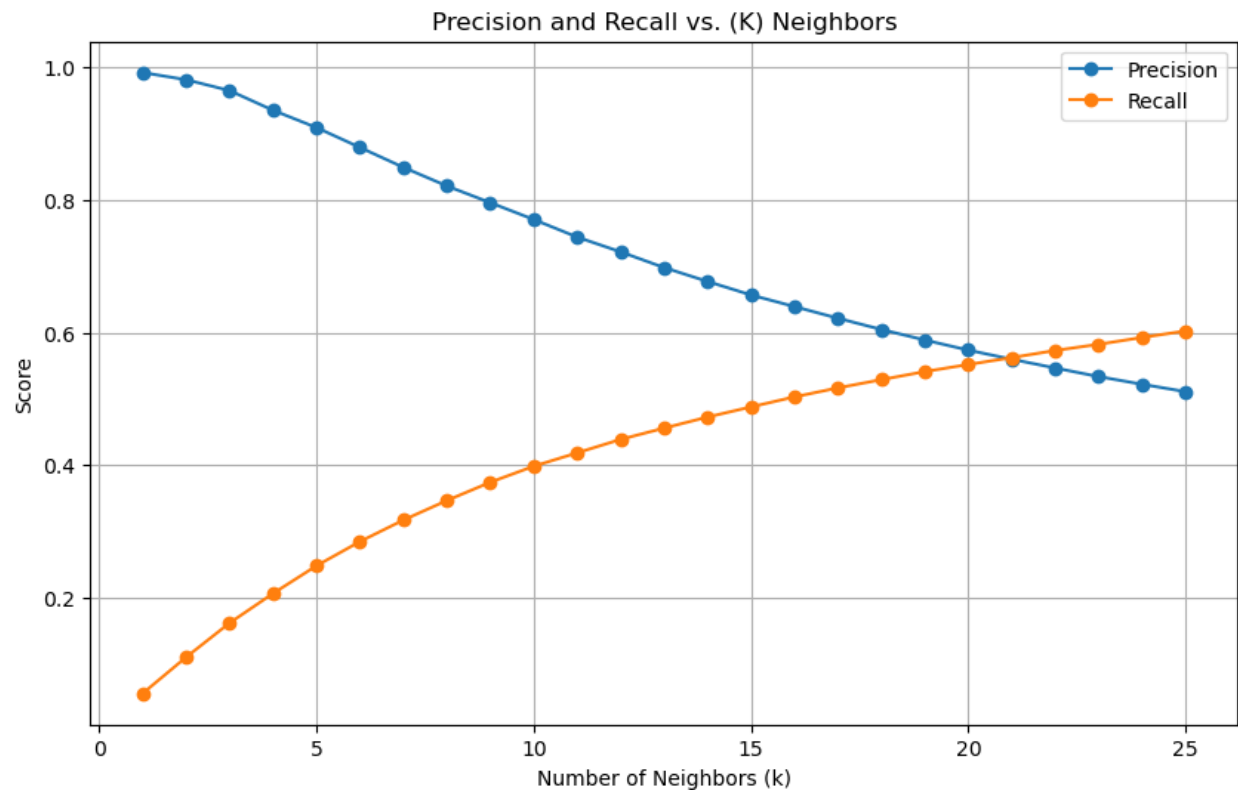
## Validation

To evaluate the performance of the k-nearest neighbors model, we employed precision and recall metrics. The precision metric evaluated how many of the selected items in the system were relevant, while the recall metric evaluated how many relevant items were selected. To calculate precision and recall, we created a matrix of predicted ratings containing the predicted ratings for all books based on their similarity to other books in the dataset. We measured the precision and recall for k values 1 through 25 in order to measure the optimal k value. For each user in the matrix, the books they rated were identified and stored. The predicted ratings for individual users were retrieved, and the top k recommendations were selected. A set intersection was performed on the retrieved top k recommendations and the actual books that the user had interacted with, the length of which was stored as the number of true positives. We then calculated the precision and recall using these equations:
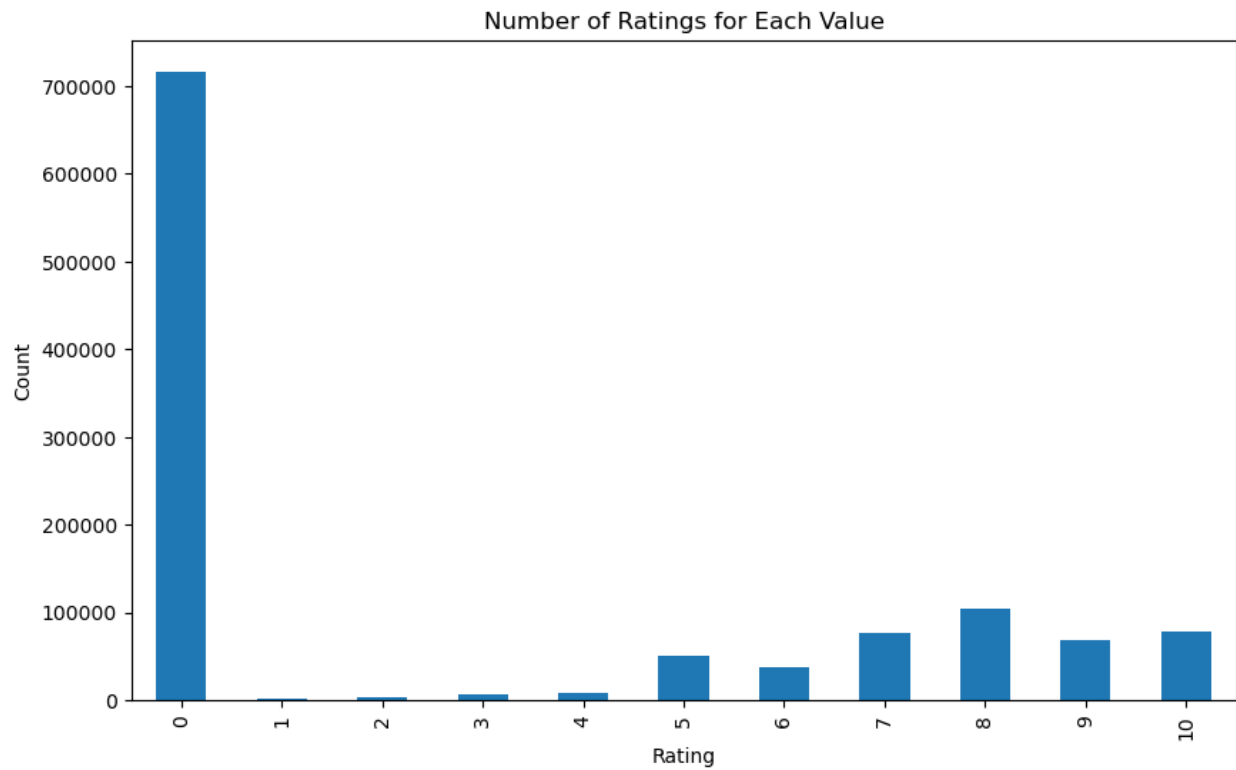
$$Precision = \frac{True\ Positives}{k}$$

$$Recall = \frac{True\ Positives}{Relevant\ Items}$$

After the calculation of precision and recall for each user, the mean value was stored for each k value. This gave us a value that was appropriate for evaluating the overall performance across all users in the matrix. This evaluation method allowed us to determine the average performance of the recommendation system, providing a metric by which we can optimize the system to enhance its performance.
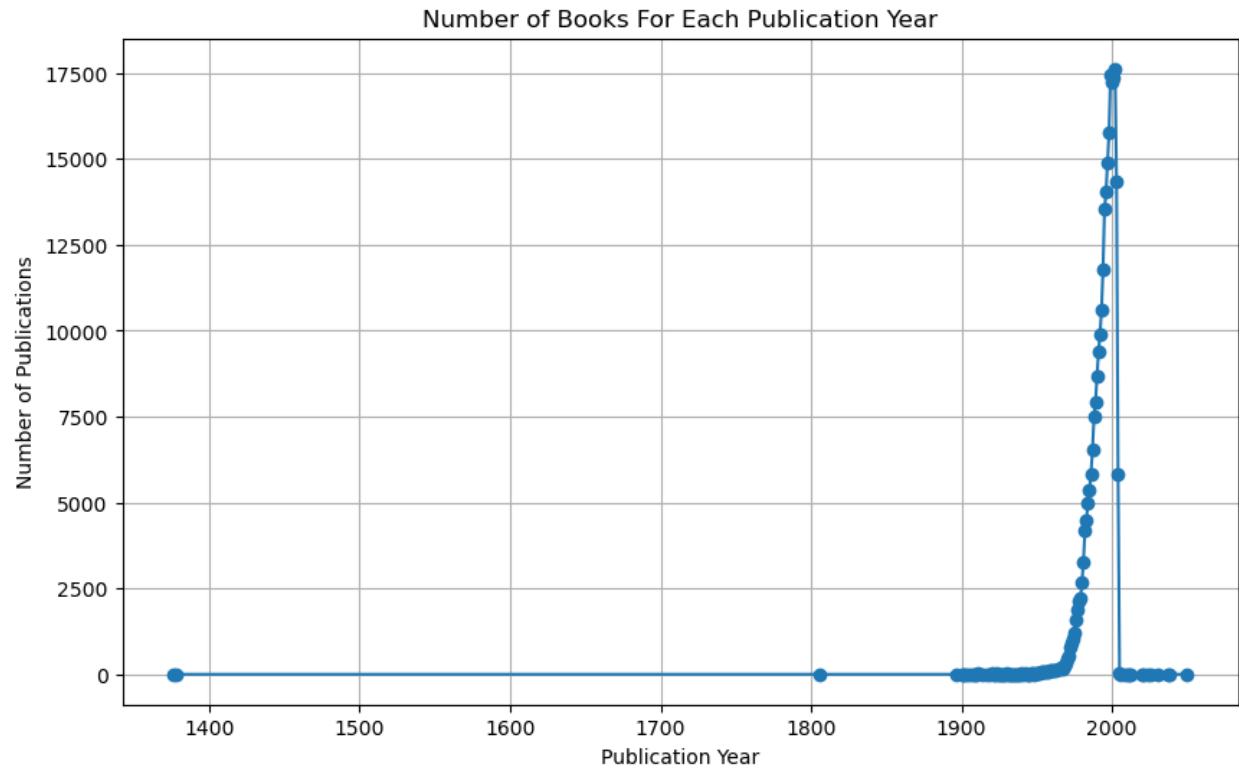
The results acquired from this validation method suggested that the optimal k value was between 10 and 13. This range represented the best balance between precision and recall, with the precision value lying between 0.77 and 0.7 and the recall value lying between 0.4 and 0.46. This suggested that the recommendations were accurate, but the model was only obtaining around 40% of possible relevant books for users. Given that we intended for the application to provide only ten recommendations based on user input, we considered the lower coverage acceptable, as the precision metric was satisfactory. A chart was created to plot the precision score against the recall score based on the number of k neighbors and help visualize the precision and recall metrics. This chart can be found below.
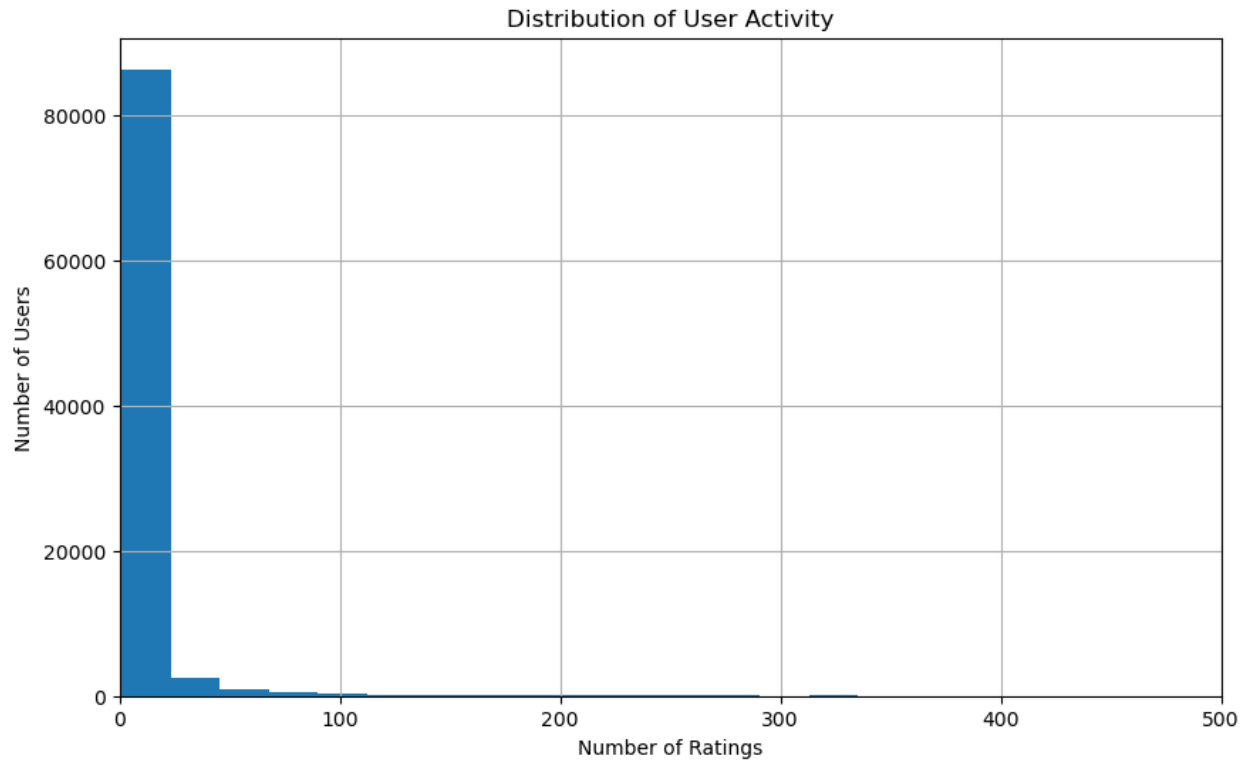
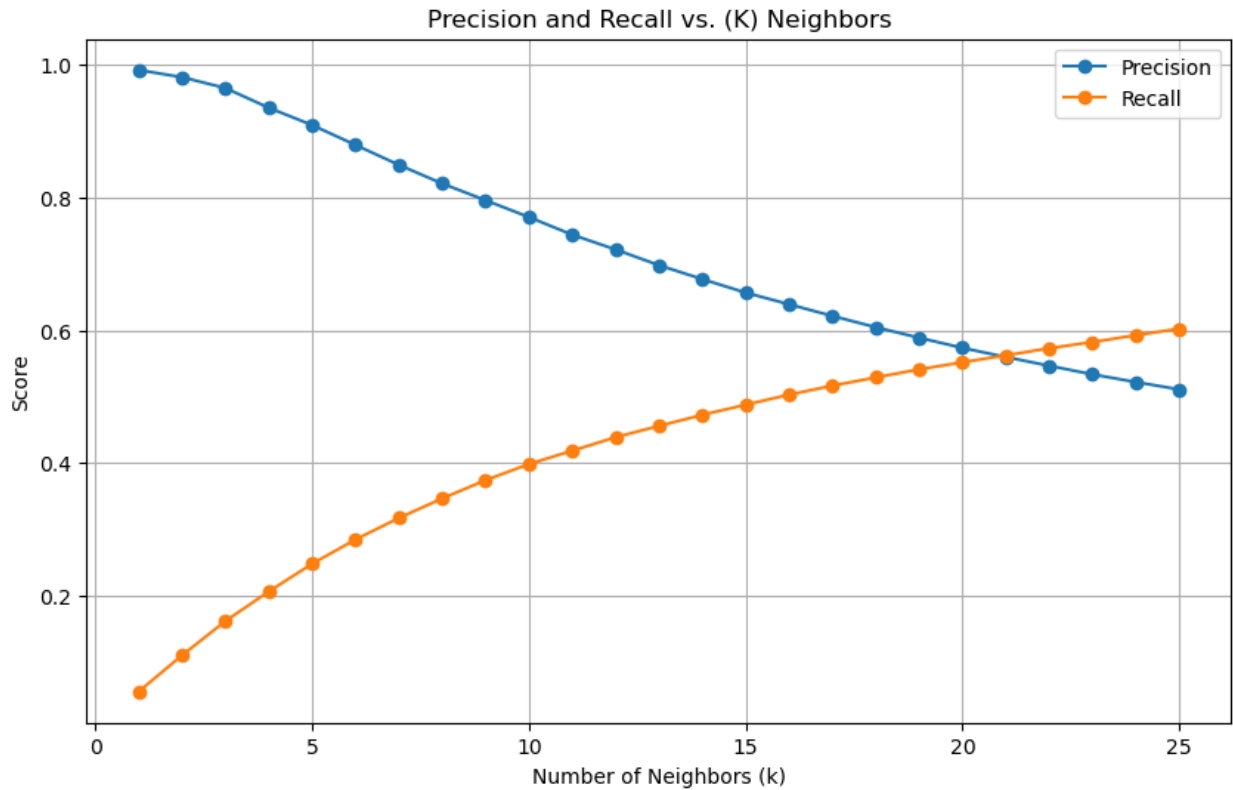Precision and Recall vs. (K) Neighbors

# Visualizations



This visualization displays the number of ratings given for every rating option (1-10) and provides insight into user behavior. This chart can be found at cell 2 in book_recommender.ipynb.
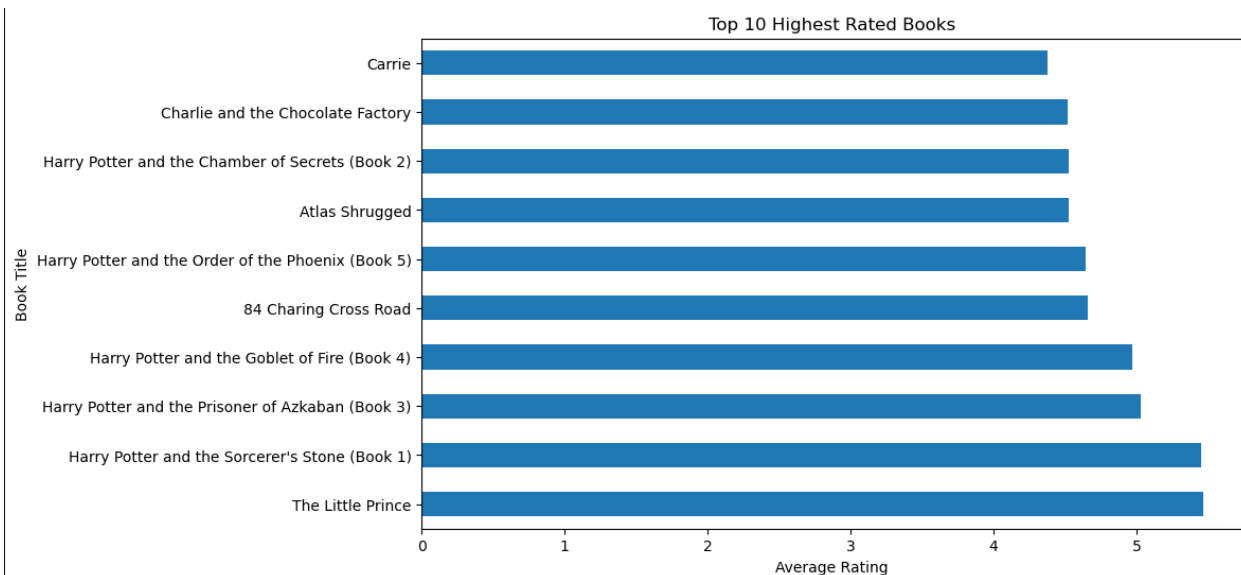
Number of Books For Each Publication Year

This visualization displays the number of books that were published each year. This chart provides information on the books that populate the dataset, as we can infer that most are from before 2000. This chart can be found at cell 3 in book_recommender.ipynb.

Distribution of User Activity

This visualization displays the number of users who have given x ratings. We inferred from this chart that the over 80,000 users who have provided fewer than 100 ratings would need to be filtered from the dataset to maintain the machine learning model's accuracy. This chart can be found at cell 5 in book_recommender.ipynb.

Precision and Recall vs. (K) Neighbors

This visualization charted the precision and recall scores against each other, allowing us to find the optimal number of k neighbors. This chart can be found at cell 11 in book_recommender.ipynb.
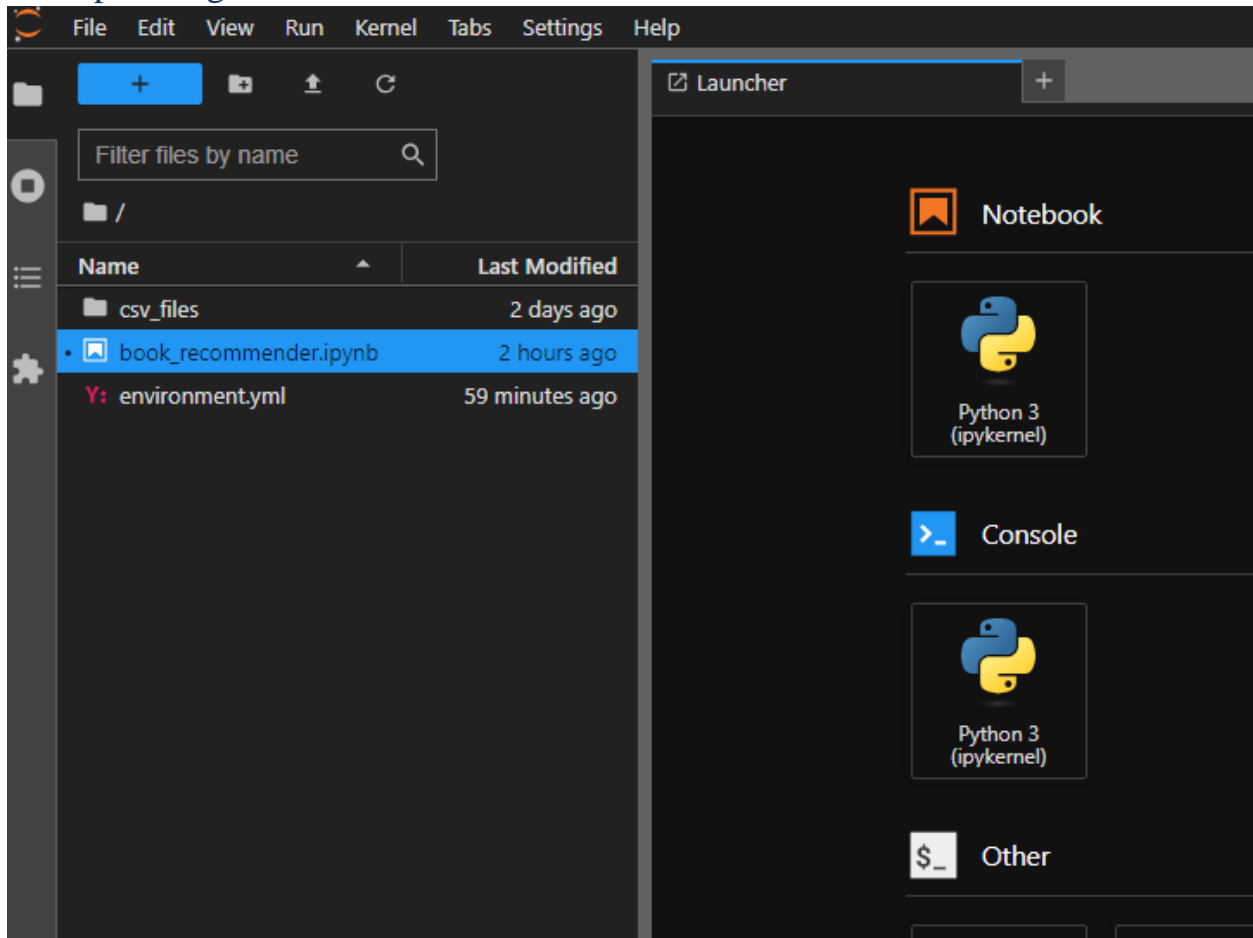


Top 10 Highest Rated Books

This chart displays the ten books with the highest ratings based on average ratings. It was used to gather book titles to use as input for testing the system. It can be found at cell 12 in book_recommender.ipynb
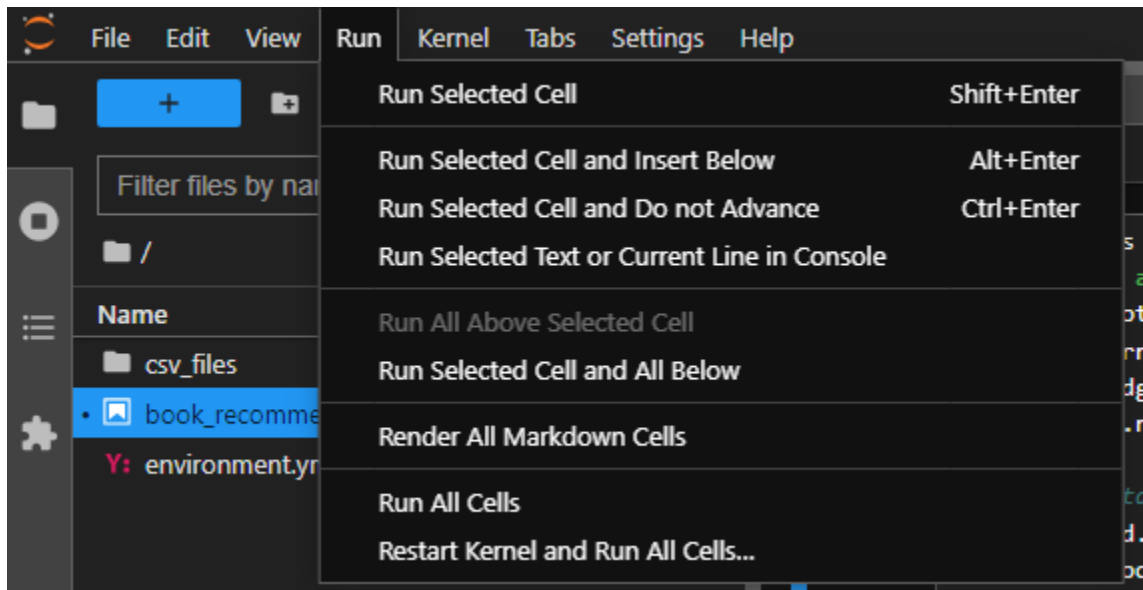
# User Guide

## How to Install

1. Open your web browser and navigate to https://docs.anaconda.com/miniconda/

2. Download the Miniconda installer for Windows 64-bit

3. Run the installer

4. Agree to the license agreement and select "Just Me" as the installation type

5. Select an install path or leave it as default, then click next

6. Leave the advanced installation options as they are, and click install

7. Once the installation is complete, you may exit the installer

8. Open the RecommenderAppCapstone.zip folder and place the recommender_application_files in "C:\Users\username\miniconda3\envs"

9. Open the Anaconda Prompt (Miniconda)

10. In the prompt, enter:

    cd C:\Users\username\miniconda3\envs\recommender_application_files

11. In the prompt, enter:

    conda env create -f environment.yml

12. After the environment has been created and the dependencies have been installed, activate the environment by typing:

    conda activate book_recommender

13. With the environment activated, type:

    jupyter lab

14. With Jupyter Lab open in your browser, double-click the book_recommender.ipynb on the file browser to access the application

15. In the top ribbon, open the "Run" tab, then click "Run All Cells"

16. When all the cells have been executed, you can access the user interface underneath the final cell (cell 14). Visualizations can be found at the locations described on pages 20-23

17. To exit the application, open the "File" tab in the top ribbon, then click "Shut Down"

## Example Usage



Double click book_recommender.ipynb.
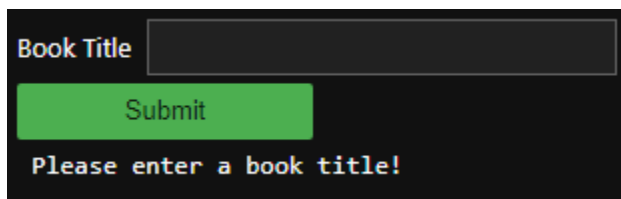


Open the "Run" tab and click "Run All Cells".

```
Book Title   The Great Gatsby

        Submit

All Creatures Great and Small
How to Be Good
The Catcher in the Rye
The Handmaid's Tale : A Novel
Falling Angels
Of Mice and Men (Penguin Great Books of the 20th Century)
Sense and Sensibility
The House on Mango Street (Vintage Contemporaries)
Fried Green Tomatoes at the Whistle Stop Cafe
Tis: A Memoir
```
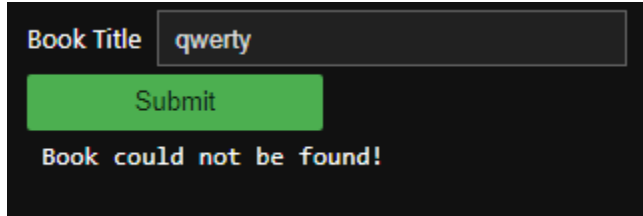
The user interface can be accessed underneath the final cell.

```
Book Title   

        Submit

Please enter a book title!
```
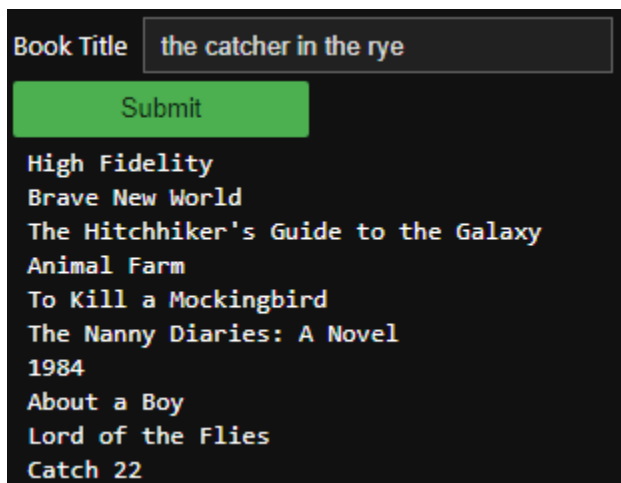
Upon submission of an empty string, this error is displayed.

```
Book Title   qwerty

        Submit

Book could not be found!
```

This error is displayed when a book with a matching title cannot be found.

```
Book Title   the catcher in the rye

        Submit

High Fidelity
Brave New World
The Hitchhiker's Guide to the Galaxy
Animal Farm
To Kill a Mockingbird
The Nanny Diaries: A Novel
1984
About a Boy
Lord of the Flies
Catch 22
```

The title may be case-insensitive.