

# AI Cap. Final Project

Team 2, Member 黃芷柔, 李宜蓁, 劉詠

## Goal

In this project, we implement a book recommendation system using two methods: user-based recommendation and content-based recommendation. We will evaluate and compare the effectiveness of each method to determine which one performs better.

## Motivation

Nowadays, many social media and shopping platforms use recommendation systems. We are undertaking this project due to our curiosity about recommendation algorithms. Through this project, we aim to deepen our understanding and become more familiar with the workings of recommendation systems.

## Dataset

Dataset: [goodbooks\\_10k\\_rating\\_and\\_description](#)

The dataset is about books with user ratings and book descriptions and is sourced from the Kaggle. Following is the introduction of the columns we will use in this project:

### **goodbooks\_10k\_rating\_and\_description.csv**

- book\_id: the index of each book
- book\_title: the title of the book
- book\_desc: the description of the book

### **ratings.csv**

- user\_id: the index of each user
- book\_id: the book the user has rated
- ratings: the ratings user gives to the books

## Methods and Techniques

### User-based

In this method, we will implement user based collaborative filtering with cosine similarity.

First, we will need to create a user-item matrix, where each row represents the users, and each column represents the books. However, we

can't directly create such a matrix from the dataset, since the matrix will be too great and our cpu can't afford it. Therefore, we will need to preprocess our data at the beginning.

Since in the dataset "rating", each user has the rating of the books that he/she has read, and in the dataset "goodbooks\_10k\_rating\_and\_description", we have more important features that we need.

Combine the features in datasets above to get the data frame with user id, book id, users' ratings of the books, books' titles. And discard the books with fewer comments, later create another data frame, "user\_book\_df", with index of user id and columns of books title, and the values in each column are the ratings that users gave that book.

Next we select a random user who is going to be recommended books for, first find the user's read books. Split those read books into two parts in if we need do evaluation later. Use the titles of the read books to extract the ratings of the books of the data frame "user\_book\_df", and save the result in "books\_read\_df". Count the numbers without null value for each user in "books\_read\_df".

The numbers count in the above represents the number of books read by each user that are the same as the random user. We define a particular ratio and percentage, if the number of a user is higher than that percentage, the user will be considered the one that has similar behavior as the random selected user.

Since we have known similar behavior users, withdraw all the rating data of those users, use the Cosine Similarity to compute similarity between those users based on their rating values of "books\_read\_df". Then we will obtain the users that are correlated with the selected random user based on the Cosine Similarity score, note that the higher the better.

For those filtered users, we further use the score obtained by the Cosine Similarity to multiply the ratings of all the books that were read by those filtered users. Therefore, each book will have many values that come from different users. We calculate the mean of those values as the weighted rating of that particular book.

Finally we choose a threshold and use the weighted rating obtained to select the books with value that surpass that threshold as recommended books for the selected random user.

## Content-based

We analyzed whether the content of the books are similar by checking the column "book\_desc\_tags\_FE", which is an introduction of every book in nearly 300 words. Then we use TF-IDF(term frequency - inverse document frequency) to vectorize the description of all books. TD-IDF is a measure of the importance

of a word, which is dependent on TD and IDF. TF is the frequency of a term in the document. If the term appears in the document more frequently, TD becomes larger, that means the word is more important. IDF is the frequency of this term in all documents. If the term appears in many documents, that means this term may be a commonly used word (such as “I”, “the”), so IDF becomes smaller. We can get TD-IDF from this equation.

$$tf(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}},$$

$f_{t,d}$ : the raw count of a term in a document

$$idf(t, D) = \log \frac{N}{|\{d : d \in D \text{ and } t \in d\}|}$$

$N$ : total number of documents in the corpus  $N = |D|$

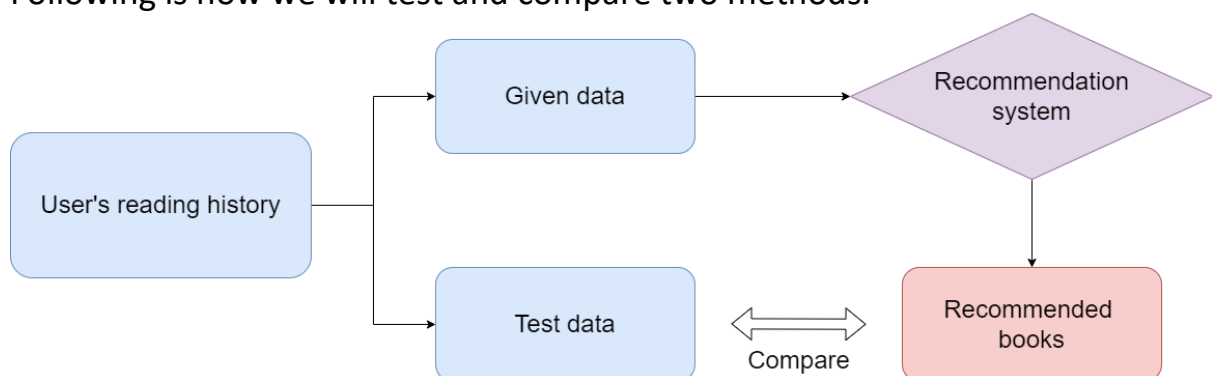
$|\{d : d \in D \text{ and } t \in d\}|$ : number of documents where the term  $t$  appear (i.e.  $tf(t, d) \neq 0$ )

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

Then, we apply cosine similarity to produce a similarity matrix with the vectorized data. This matrix allows us to look up the relationships between books. To recommend books, we choose a book, find the data related to this book in the similarity matrix and extract them. Next, we sort the data by similarity, make the books with higher similarity in front of the list. We can get the five books closest to the selected book.

## Testing & Result

Following is how we will test and compare two methods.



We wanted to know if the books we recommend match the readers' preferences, so we analyzed our recommendation system by examining the users' read history. We chose some random users, divided the read history of each user into 2 parts. We put one part into the recommendation system to produce recommended books. Then we checked whether these books are in the

other part of the user's read history. If there are books matching, that means the user may be interested in what we recommended.

In a user-based system, we split the users' read history into two parts to put one part in the recommendation system directly. In a content-based system, since the unit of analysis is a single book, we analyze the viewing records book by book for recommendations. If there are duplicates, they are removed, and the number of recommended books is adjusted to match the quantity in the user-based system.

We let:  $Matching\ rate = \frac{The\ number\ of\ matching\ books}{Total\ test\ data\ of\ an\ user}$

while the scale of recommended books is same in two methods (almost 500 books while the number of total books is 10000)

Result of 30 random users:

	User-based	Content-based
Average of matching rate	33.28%	59.88%

We can find that the content-based model overperforms the user-based model. We speculate that this is because the preferences between users who read the same books are not that close.

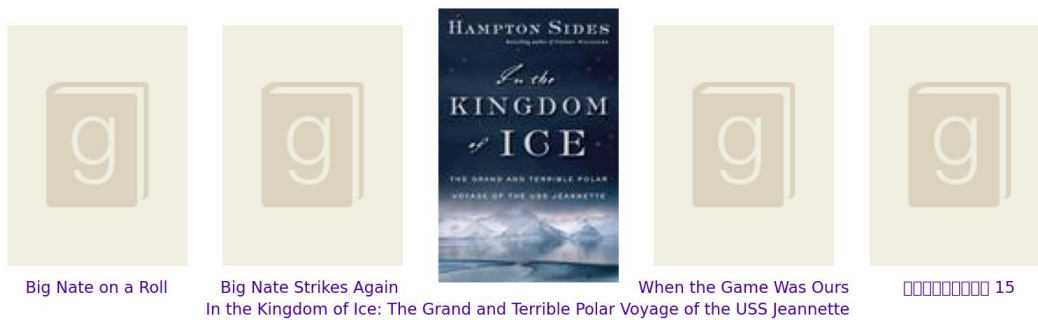
## Demo

### 1. user-based

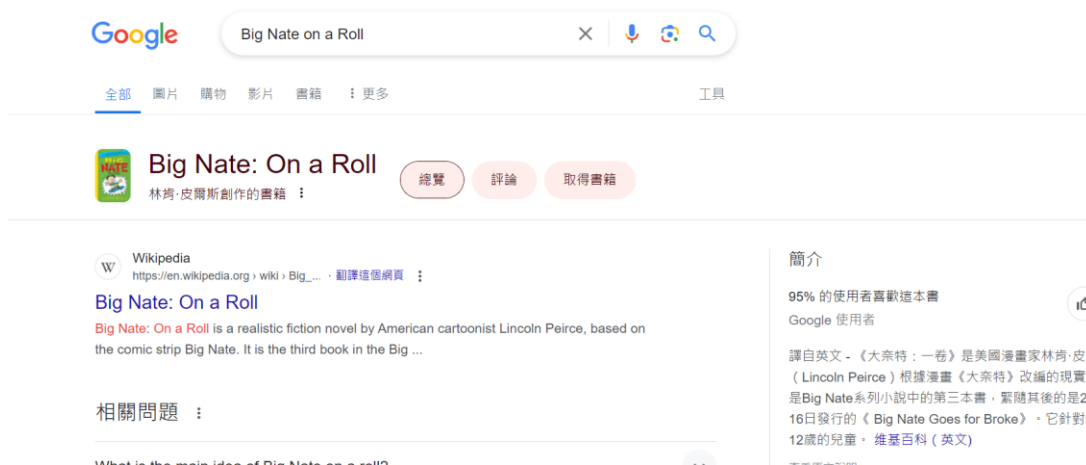
```
select user: 23132
Sample ten books the user has read:
[0] Confessions of an Ugly Stepsister
[1] Firefly Lane
[2] My Last Duchess
[3] Once a Runner
[4] The Language of Flowers
[5] The Scarlet Letter
[6] The Giving Tree
[7] Bergdorf Blondes
[8] 人生がときめく片づけの魔法
[9] Last Night at Chateau Marmont

Recommend books:
Big Nate on a Roll (https://www.google.com/search?q=Big+Nate+on+a+Roll)
Big Nate Strikes Again (https://www.google.com/search?q=Big+Nate+Strikes+Again+)
In the Kingdom of Ice: The Grand and Terrible Polar Voyage of the USS Jeannette (https://www.google.com/search?q=In+the+Kingdom+of+Ice%3A+The+Grand+and+Terrible+Polar+Voyage+of+the+USS+Jeannette)
When the Game Was Ours (https://www.google.com/search?q=When+the+Game+Was+Ours)
フルーツバスケット 15 (https://www.google.com/search?q=%E3%83%95%E3%83%AB%E3%83%BC%E3%83%84%E3%83%90%E3%82%B9%E3%82%B1%E3%83%83%E3%83%88+15)
Mistakes Were Made (But Not by Me): Why We Justify Foolish Beliefs, Bad Decisions, and Hurtful Acts (https://www.google.com/search?q=Mistakes+Were+Made+But+Not+by+Me%29%3A+Why+We+Justify+Foolish+Beliefs%2C+Bad+Decisions%2C+and+Hurtful+Acts)
The Love Song of Miss Queenie Hennessy (https://www.google.com/search?q=The+Love+Song+of+Miss+Queenie+Hennessy)
The Holiness of God (https://www.google.com/search?q=The+Holiness+of+God)
A Higher Call: An Incredible True Story of Combat and Chivalry in the War-Torn Skies of World War II (https://www.google.com/search?q=A+Higher+Call%3A+An+Incredible+True+Story+of+Combat+and+Chivalry+in+the+War-Torn+Skies+of+World+War+II)
Willpower: Rediscovering the Greatest Human Strength (https://www.google.com/search?q=Willpower%3A+Rediscovering+the+Greatest+Human+Strength)
```

- First, select a random user and show ten of the books the user has read. Next, we will show the book recommended to the user.



- We will also show the book title and their corresponding title if they have.



- We will provide the url of google search page for each book. The url is printed behind the recommended book.

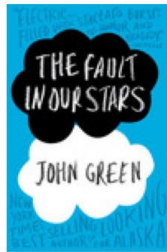
## 2. content based

```

Target book:
The Fault in Our Stars
Target book description:
Despite the tumor-shrinking medical miracle that has bought her a few years, Hazel has never been anything but terminal, her final chapter inscribed upon diagnosis. But when a gorgeous plot twist named Augustus Waters suddenly appears at Cancer Kid Support Group, Hazel's story is about to be completely rewritten. Insightful, bold, irreverent, and raw, The Fault in Our Stars is award-winning author John Green's most ambitious and heartbreaking work yet, brilliantly exploring the funny, thrilling, and tragic business of being alive and in love.
Recommend books:
The Darkest Part of the Forest (https://www.google.com/search?q=The+Darkest+Part+of+the+Forest)
Wise Blood (https://www.google.com/search?q=Wise+Blood)
The Emperor of All Maladies: A Biography of Cancer (https://www.google.com/search?q=The+Emperor+of+All+Maladies%3A+A+Biography+of+Cancer)
The Green Mile, Part 6: Coffey on the Mile (https://www.google.com/search?q=The+Green+Mile%2C+Part+6%3A+Coffey+on+the+Mile)
The Magician's Assistant (https://www.google.com/search?q=The+Magician%2S+Assistant)
Benediction (https://www.google.com/search?q=Benediction)
The House of Hades (https://www.google.com/search?q=The+House+of+Hades)
Autobiography of a Face (https://www.google.com/search?q=Autobiography+of+a+Face)
Second Chance (https://www.google.com/search?q=Second+Chance)
The E-Myth Revisited: Why Most Small Businesses Don't Work and What to Do About It (https://www.google.com/search?q=The+E-Myth+Revisited%3A+Why+Most+Small+Businesses+Don%27t+Work+and+What+to+Do+About+It)

```

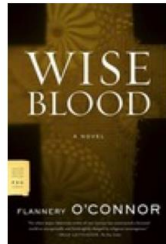
- First we will choose a target book, and print its description. Next, we will show the recommended book.



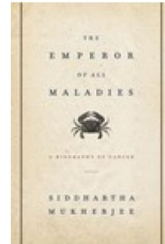
Rating: 4.2



Rating: 3.9



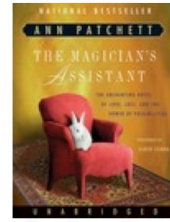
Rating: 3.9



Rating: 4.3



Rating: 4.6



Rating: 3.6

- Here is the cover of the recommended books and their ratings.



The Darkest Part of the Forest



全部 購物 圖片 影片 新聞 更多

工具

提示：顯示繁體中文的搜尋結果。你也可以進一步瞭解如何依語言篩選結果。



## The Darkest Part of the Forest

荷莉·布萊克創作的書籍

總覽

影片

評論

取得書籍



Goodreads

<https://www.goodreads.com/show> · 翻譯這個網頁

### The Darkest Part of the Forest by Holly Black

2015年1月13日 — Hazel lives with her brother, Ben, in the strange town of Fairfold where humans and fae exist side by side. The faeries' seemingly harmless ...

★★★★★ 評分：3.9 · 9,856 則評論 · US\$11.99

#### 簡介

3.9/5 · Goodreads

95% 的使用者喜歡這本

Google 使用者

譯自英文 林林墨墨

- We will provide the url of google search page for each book. The url is printed behind the recommended book.

## Conclusion

From the results, we can see that our recommendation systems can roughly infer the readers' preferences. Since the datasets we have do not contain recently user records, however, if we can solve the problem and use this system in the place such as library, or online books stores, and so on, we think it is a good way to let the users know the potential books that he or she may want to read, moreover, we have two ways to to recommendation, so we hope using these two methods, we can further explore users' preference. The above is the whole of our final project report, thank you.

# Reference

[1]Content-Based

<https://www.kaggle.com/code/mechatronixs/recommend-project-content-based-recommendation>

[2]User-Based

<https://www.kaggle.com/code/mechatronixs/recommendation-project-3-user-based-recommendation/notebook>

[3]TF-IDF Reference

<https://zh.wikipedia.org/zh-tw/Tf-idf>

[4]Cosine Similarity

[https://en.wikipedia.org/wiki/Cosine\\_similarity](https://en.wikipedia.org/wiki/Cosine_similarity)

[5]Collaborating Filtering

<https://realpython.com/build-recommendation-engine-collaborative-filtering/>

# Contributions

黃芷柔: content-based, evaluation

李宜蓁: find information, demo code

劉詠: user-based