

Book Recommendation System

Team Member



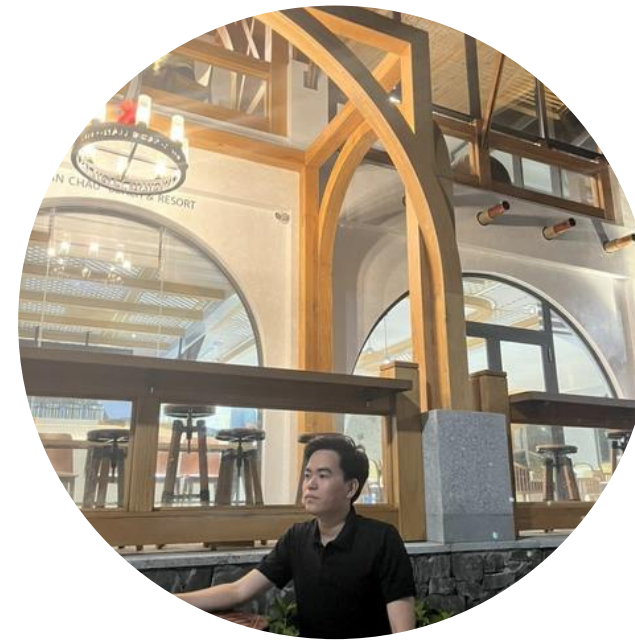
Lê Hoàng
Khang

- Write Proposal
- Search and collect data
- Content-based filtering algorithm



Nguyễn Thanh Hùng

- Writing Presentations
- Data preprocessing
- Collaborative filtering algorithm



Nguyễn Duy
Thái

- Write Report
- Search and collect data
- Content-based filtering algorithm



Lê Phước
Yên

- Write Milestone
- EDA
- Collaborative filtering algorithm

Table Of Content

01

REASON FOR
CHOOSING TOPIC

02

DATASET

03

EDA

MODEL
DEVELOPMENT

04

EXPERIMENTS AND
EVALUATION

05

CONCLUSION

06



REASON FOR CHOOSING TOPIC

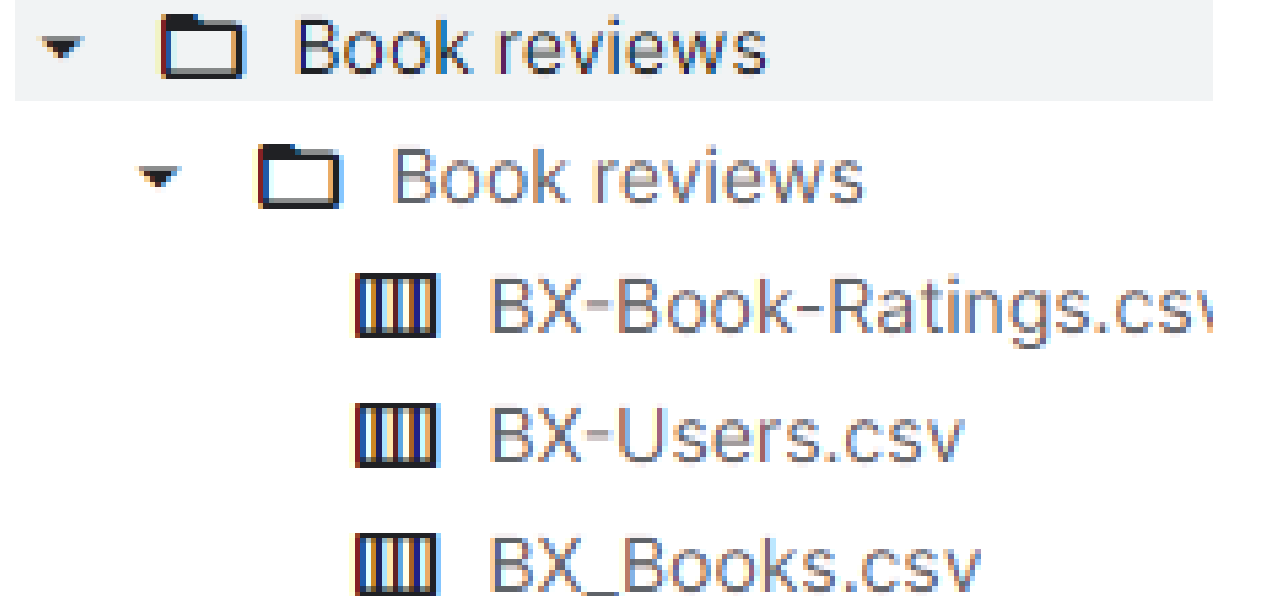


ABOUT DATASET

- Link: <https://www.kaggle.com/datasets/ruchi798/bookcrossing-dataset>
- The data set consists of 3 main files containing about 278,858 users, providing more than 1 million ratings on about 270,379 books:
- BX-Book-Ratings: includes the user's ID, the code of each book, and the user's rating for that book.
- BX-Users: includes user ID, address and age.
- BX-Books: includes each book's code, book title, author, year of publication, publisher, summary and link to the book's photo.

Data Explorer

Version 3 (600.34 MB)



EDA

Join 3 tables

```
1 df = ratings.join(users, on='User-ID')
2 df = df.join(books, on='ISBN')
✓ 0.0s
```

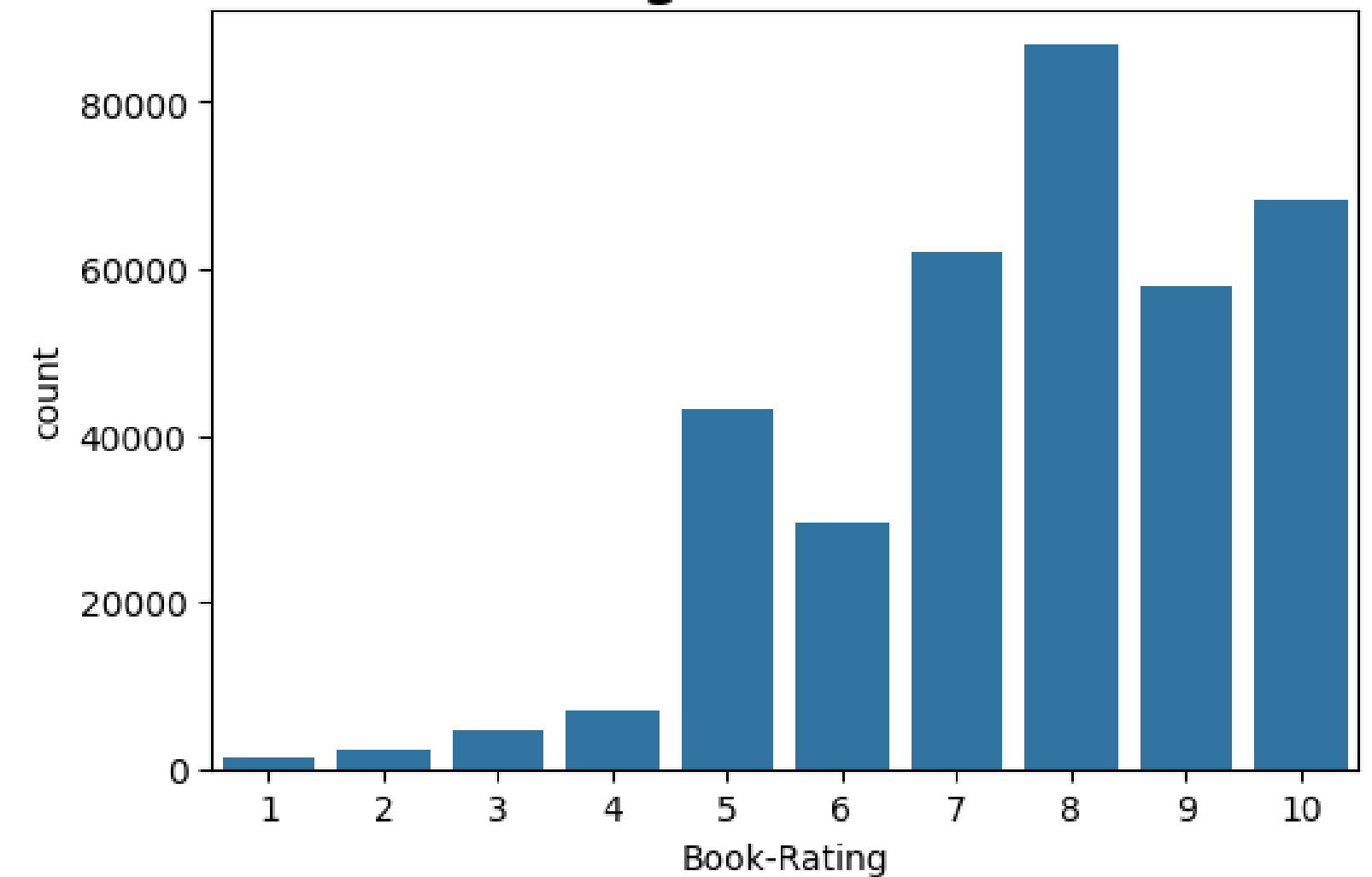
```
1 df.count()
✓ 4.7s
```

632238

Convert data types

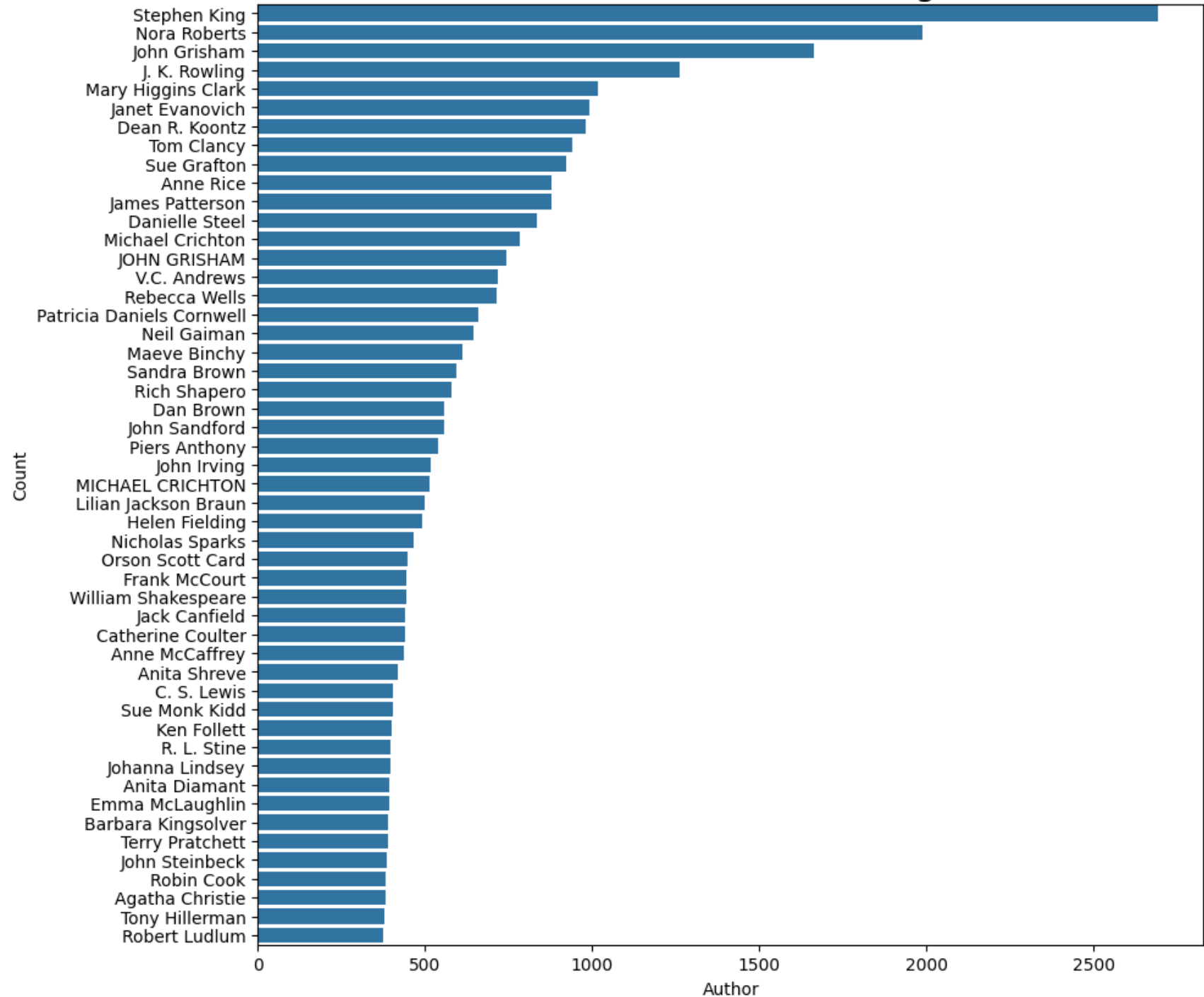
```
df = df.withColumn("Age", df["Age"].cast(IntegerType()))
df = df.withColumn("Year-Of-Publication",
                  df["Year-Of-Publication"].cast(IntegerType()))
df = df.withColumn("Book-Rating", df["Book-Rating"].cast(IntegerType()))
df = df.withColumn("User-ID", df["User-ID"].cast(IntegerType()))
```

Rating Distribution

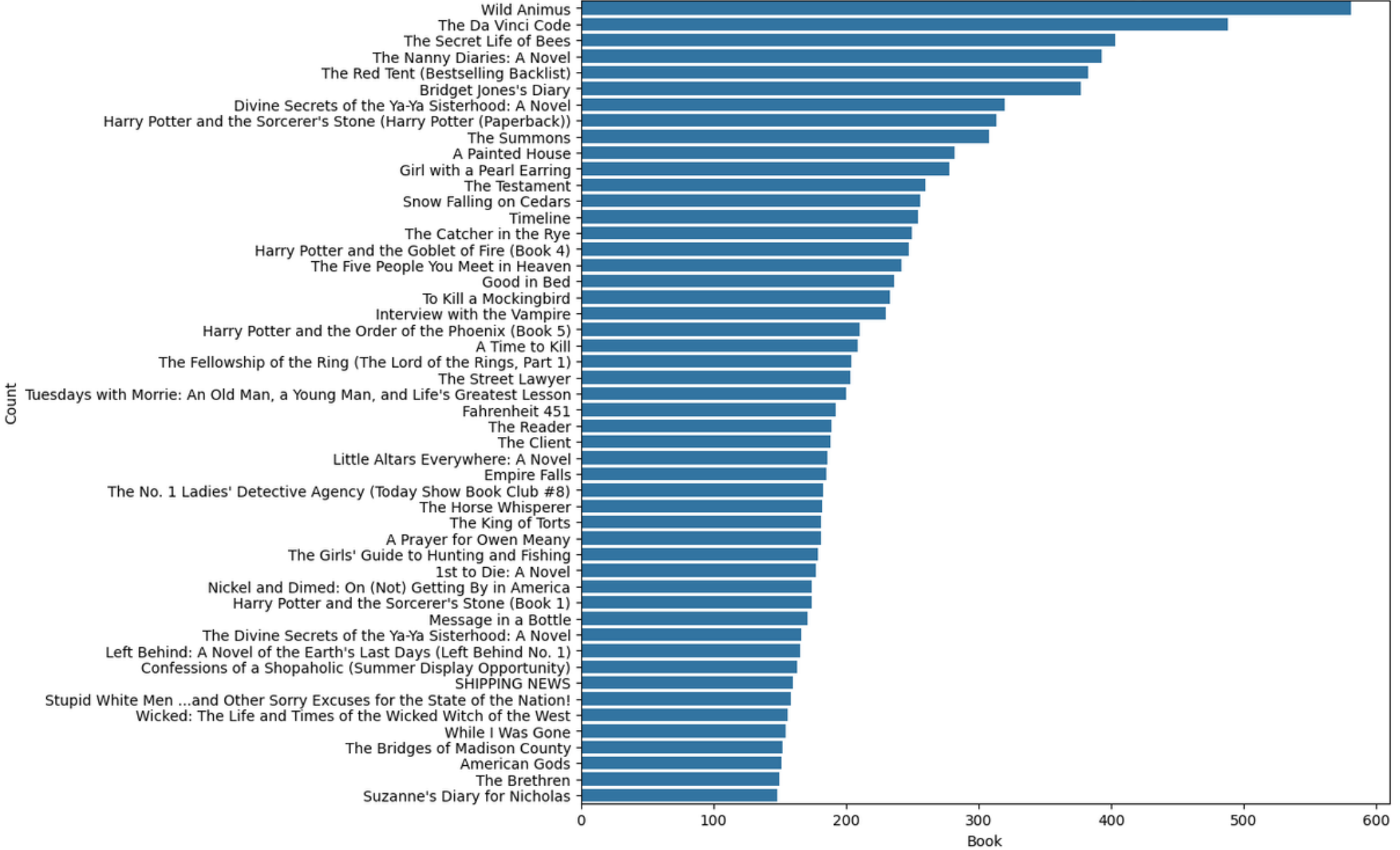




Authors with most Ratings

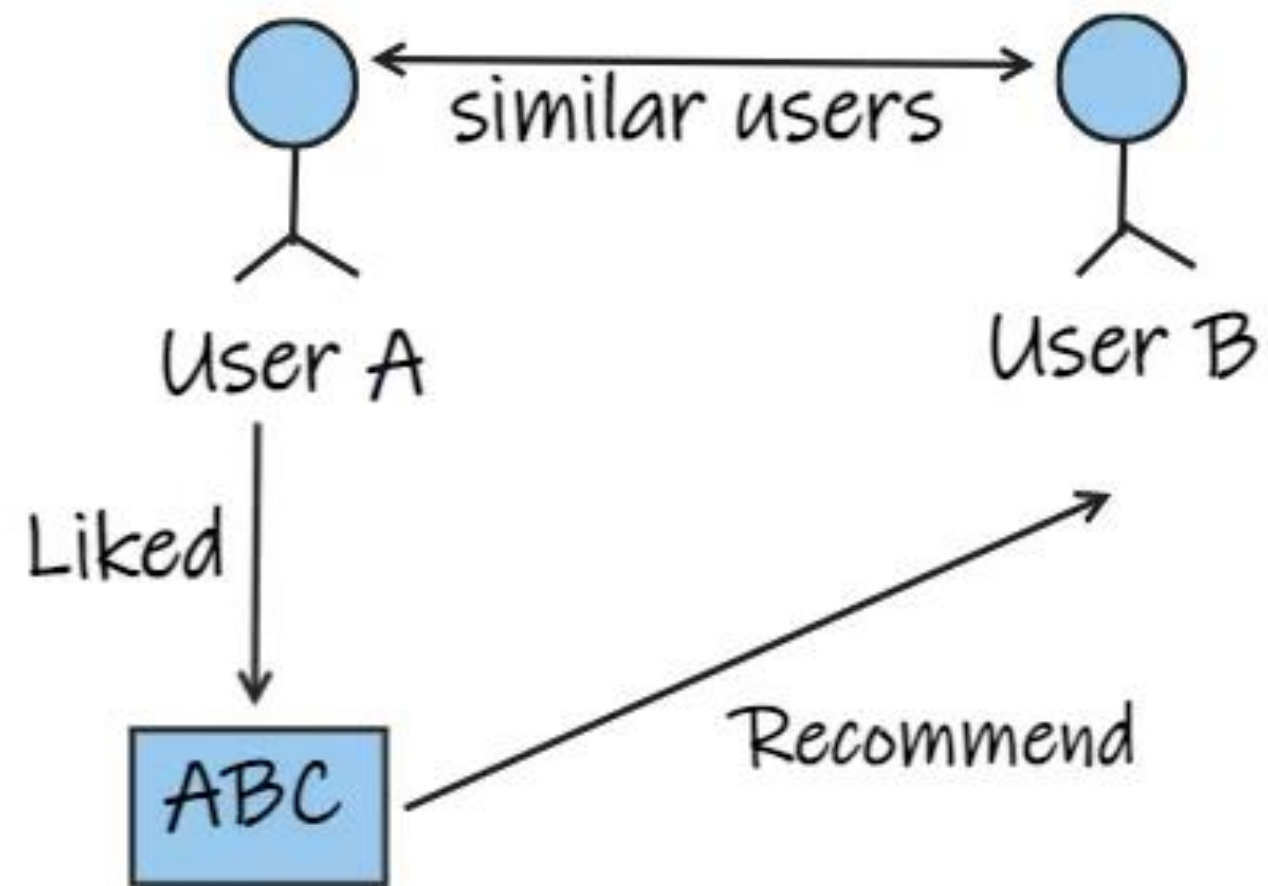


Books with most Ratings

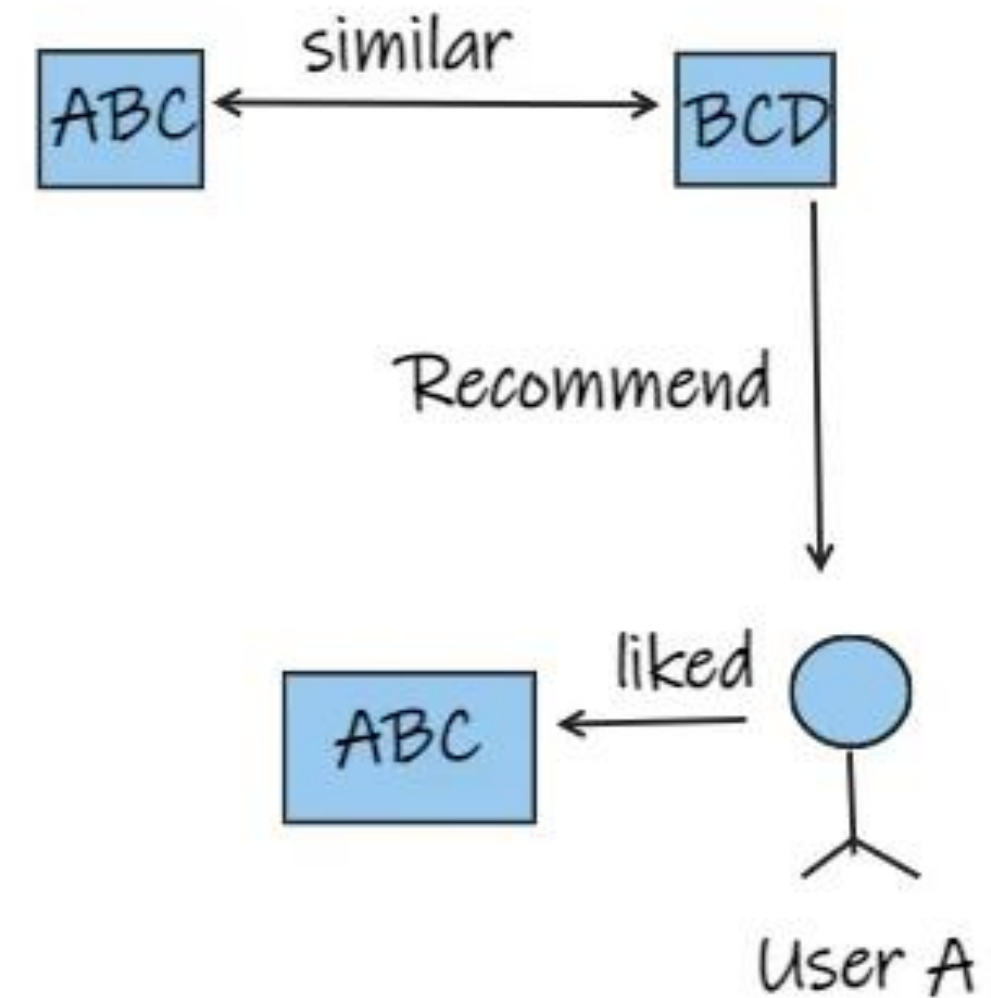


MODEL DEVELOPMENT

Collaborative Filtering

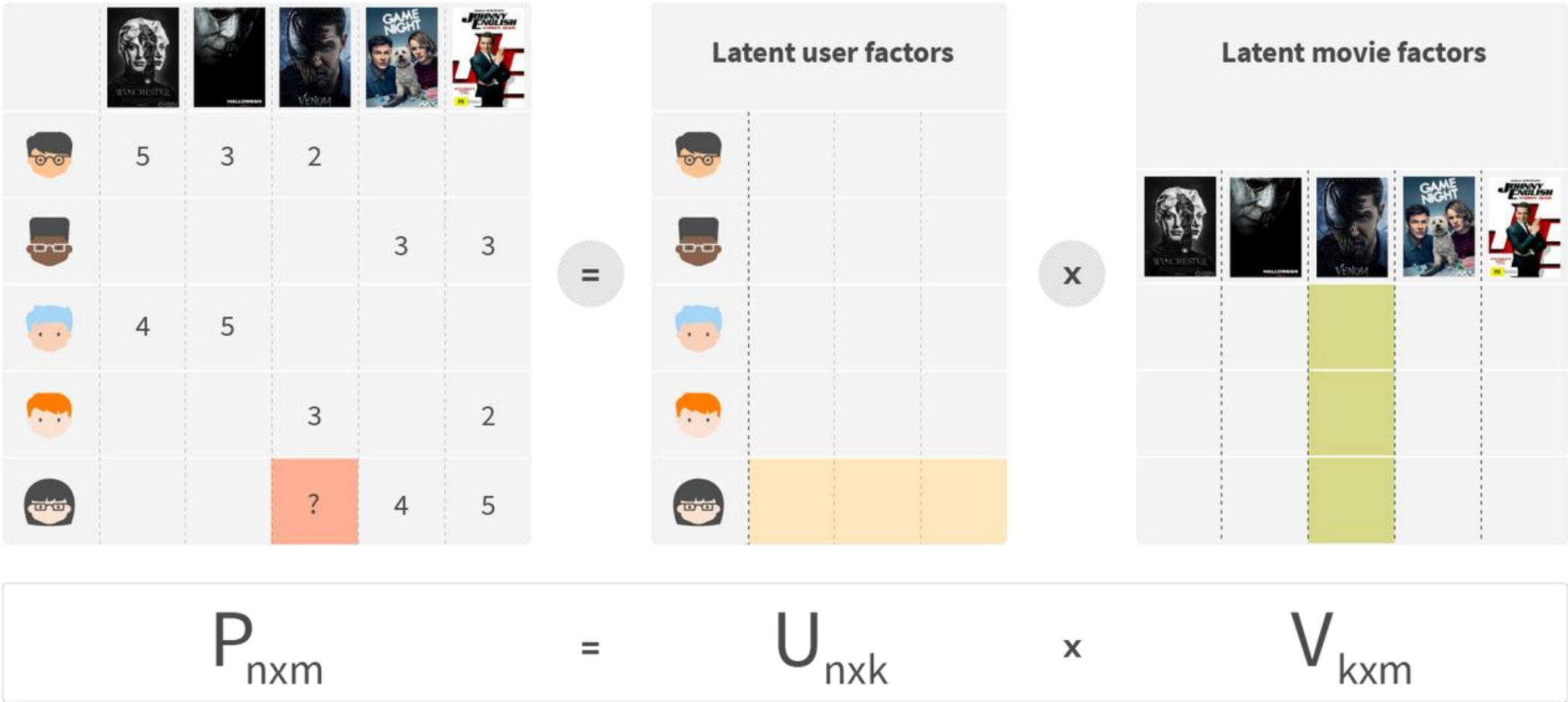


Content Based



MODEL DEVELOPMENT

COLLABORATIVE FILTERING USING ALS



Parameter	Description	Default value
rank	Số lượng latent factors	10
regParam	Regularization parameter	1
maxIters	Số lần lặp tối đa	10

MODEL DEVELOPMENT

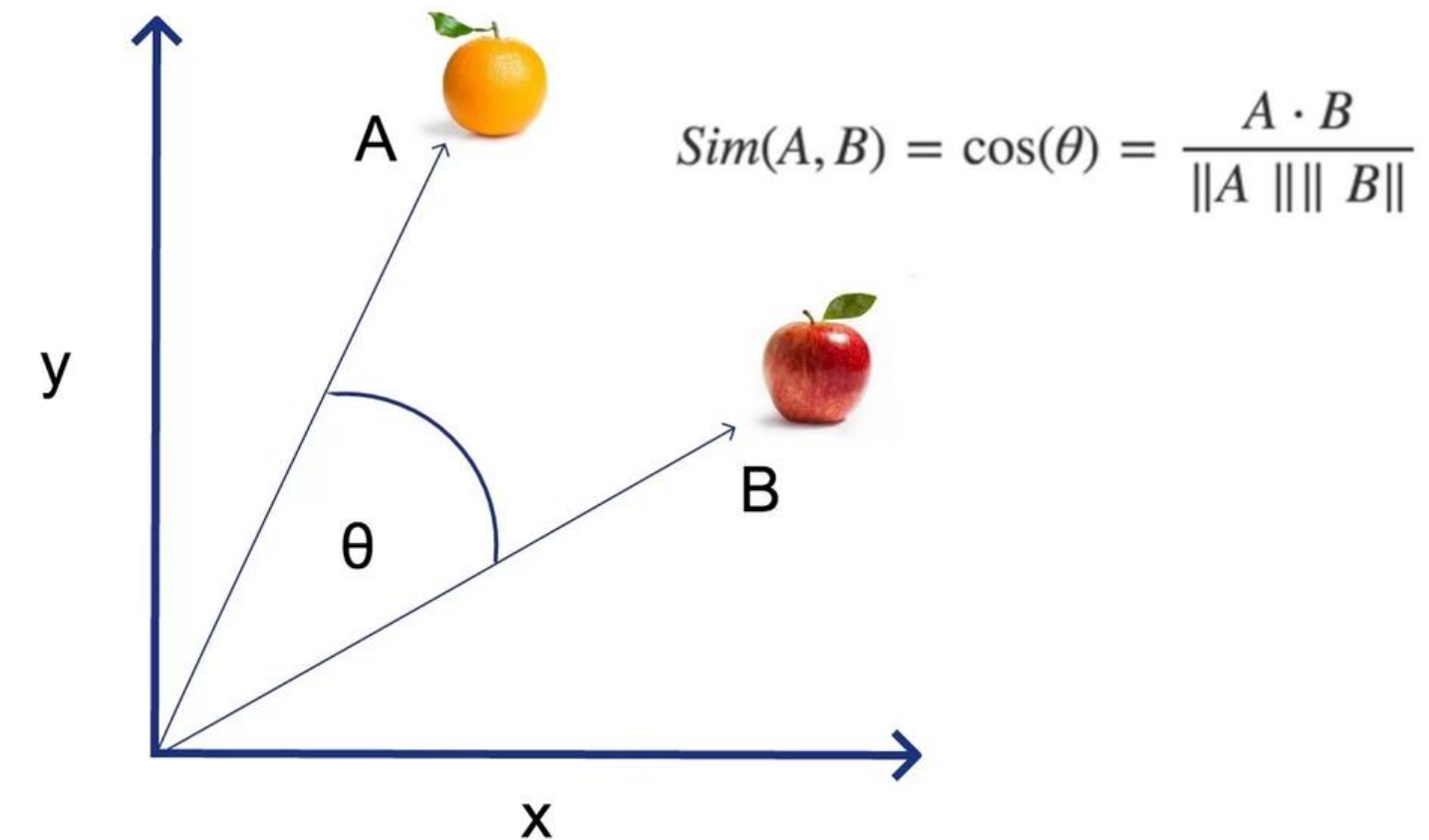
CONTENT-BASED USING TF-IDF AND COSIN SIMILARITY

$$\text{Term Frequency } TF(t, d) = \frac{\text{count of } t \text{ in } d}{\text{number of words in } d}$$

$$\text{Inverse Document Frequency} = \log_2\left(\frac{\text{number of documents } N}{\text{number of documents containing the term } t}\right)$$

$$tf - idf(d, t) = tf(t, d) * idf(t)$$

Cosine Similarity



Experiment and Evaluate

Collaborative Filtering

```
1 for_one_user = predictions.filter(col("User-ID") == 98391).select(  
2     "User-ID", "Book-Title", "Image-URL-M", "Book-Rating", "prediction"  
3 )  
4 for_one_user.show(5, False)
```

[42]

Python

```
... +-----+-----+-----+-----+-----+-----+  
|User-ID|Book-Title|Image-URL-M|Book-Rating|prediction|  
+-----+-----+-----+-----+-----+-----+  
|98391|Flirting With Trouble (Harlequin American Romance Series)|http://images.amazon.com/images/P/037375020X.01.MZZZZZZZ.jpg|8|3.8005092|  
|98391|Anything Goes (Harlequin Blaze, 112)|http://images.amazon.com/images/P/037379116X.01.MZZZZZZZ.jpg|9|7.2175493|  
|98391|The Rake: Lessons in Love|http://images.amazon.com/images/P/038082082X.01.MZZZZZZZ.jpg|8|7.8791027|  
|98391|Key of Light (Key Trilogy (Paperback))|http://images.amazon.com/images/P/051513628X.01.MZZZZZZZ.jpg|10|7.6027946|  
|98391|The Reluctant Suitor|http://images.amazon.com/images/P/0060185708.01.MZZZZZZZ.jpg|9|8.014934|  
+-----+-----+-----+-----+-----+-----+
```

only showing top 5 rows

RMSE value 2.21967872716262

Experiment and Evaluate

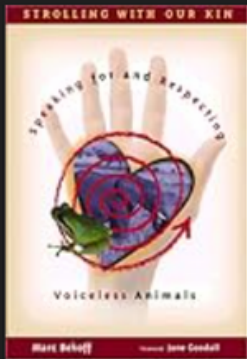
Content-Based Filtering

Books similar to keyword - wild animal

ISBN	similarity_score	Book-Title	Book-Author	Year-Of-Publication	Publisher
059031971X	0.6469376166083098	The Gentle Jungle	Toni R. Helfer	1981	Scholastic Pa
1564585506	0.595001952867062	Big Pictures: Wild Animals	Mary Ling	1994	Dorling Kinde
0440159415	0.5845439562587951	My Wild World	Joan Embery	1981	Bantam Books
1881699021	0.5478580272597805	Strolling with Our Kin: Speaking for and Respecting Voiceless Animals	Marc Bekoff	2000	Amer Anti-Viv.
0486296547	0.5397765730533401	Riddles, Riddles, Riddles (Dover Game & Puzzle Activity Books)	Darwin A. Hindman	1997	Dover Publica

only showing top 5 rows

Strolling with Our Kin: Speaking for and Respecting Voiceless Animals



A look at animal protection and compassion, using animals as fur and food, animal pain and suffering, and dissection and vivisection.

Big Pictures: Wild Animals



Each photograph of a wild mammal or bird is accompanied by a question and answer offering information about the animal or its behavior

The Mystery of the Ivory Charm (Nancy Drew Mystery Stories, No 13)



Nancy Drew determines whether an ivory elephant charm really protects its wearer from harm when she investigates the involvement of a member of the Bengleton Wild-Animal Show in a mysterious scheme.

A close-up photograph of a person's hand holding a silver pen, pointing at a document. The document features a network diagram with nodes and connecting lines. The background is slightly blurred, showing a laptop screen and other papers.

Conclusion and development direction

Conclusion

- Through the process of implementing the project, the group learned more knowledge about recommendation models using PySpark for book suggestions.
- Overall, the topic has completed the basic level, providing suggestions based on user preferences and entered keywords.

Development direction

- Perform suggestions with large amounts of data
- Combine interfaces to be more intuitive
- Process automation



Thanks