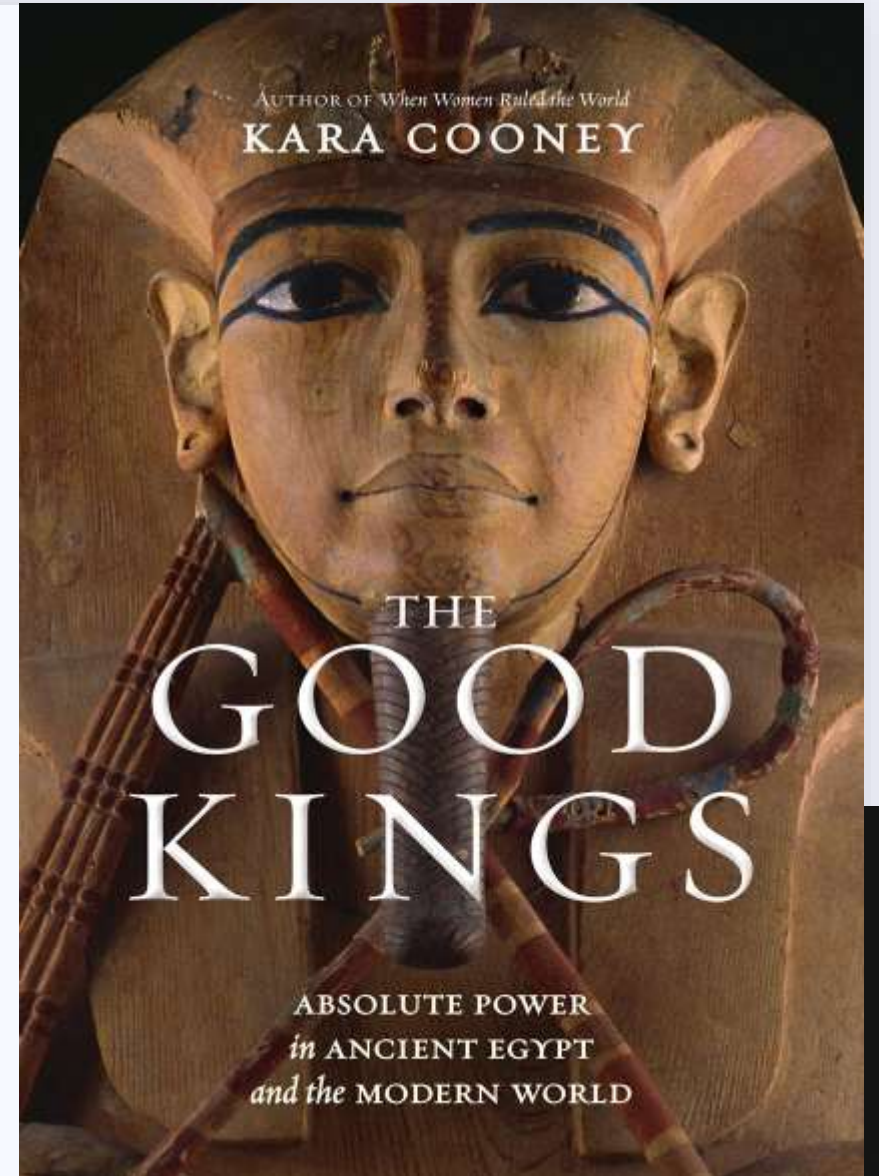# BOOK RECOMMENDER

# INTRODUCTION

The project is primarily focused on building a **book recommendation system**, which typically falls under the domain of **collaborative filtering**. Collaborative filtering is an unsupervised learning approach used for recommendation tasks. It's neither regression nor classification, as the goal is to predict user preferences or item similarities rather than a specific numeric value or category.

The **k-nearest neighbors (KNN) algorithm** is commonly used for collaborative filtering in recommendation systems. In this case, it's used to find similar books based on user preferences or book characteristics. KNN is a type of unsupervised learning algorithm that belongs to the broader category of instance-based learning.

In the context of a book recommendation system, the primary learning task is generally considered unsupervised, as the goal is to find **patterns or similarities in user behavior or item features without explicit labels**.

# THE DATASET

The dataset was obtained from Kaggle:
https://www.kaggle.com/datasets/jealousleopard/goodreadsbooks?resource=download

- 12 columns
- 11126 unique entries

# EXPLORING DATA



```
df.head()
```

| | bookID | title | authors | average_rating | isbn | isbn13 | language_code | num_pages | ratings_count | text_reviews_count | publication_date |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | Harry Potter and the Half-Blood Prince (Harry ... | J.K. Rowling/Mary GrandPré | 4.57 | 0439785960 | 9780439785969 | eng | 652 | 2095690 | 27591 | 9/16/2006 |
| 1 | 2 | Harry Potter and the Order of the Phoenix (Har... | J.K. Rowling/Mary GrandPré | 4.49 | 0439358078 | 9780439358071 | eng | 870 | 2153167 | 29221 | 9/1/2004 |
| 2 | 4 | Harry Potter and the Chamber of Secrets (Harry... | J.K. Rowling | 4.42 | 0439554896 | 9780439554893 | eng | 352 | 6333 | 244 | 11/1/2003 |

# EXPLORING DATA



```
#Exploring the dataset:
print(df.info())
print(df.describe())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11123 entries, 0 to 11122
Data columns (total 12 columns):
 #   Column              Non-Null Count  Dtype
---  ------              --------------  -----
 0   bookID              11123 non-null  int64
 1   title               11123 non-null  object
 2   authors             11123 non-null  object
 3   average_rating      11123 non-null  float64
 4   isbn                11123 non-null  object
 5   isbn13              11123 non-null  int64
 6   language_code       11123 non-null  object
 7    num_pages          11123 non-null  int64
 8   ratings_count       11123 non-null  int64
 9   text_reviews_count  11123 non-null  int64
 10  publication_date    11123 non-null  object
 11  publisher           11123 non-null  object
dtypes: float64(1), int64(5), object(6)
memory usage: 1.0+ MB
None
             bookID  average_rating        isbn13     num_pages  \
count  11123.000000    11123.000000  1.112300e+04  11123.000000
mean   21310.856963        3.934075  9.759880e+12    336.405556
std    13094.727252        0.350485  4.429758e+11    241.152626
min        1.000000        0.000000  8.987060e+09      0.000000
25%    10277.500000        3.770000  9.780345e+12    192.000000
```

# NULL VALUES



```
6]:  #finding null values:
     df.isnull().sum()

6]:  bookID                 0
     title                  0
     authors                0
     average_rating         0
     isbn                   0
     isbn13                 0
     language_code          0
        num_pages           0
     ratings_count          0
     text_reviews_count     0
     publication_date       0
     publisher              0
     dtype: int64
```

# DATA VISUALIZATION

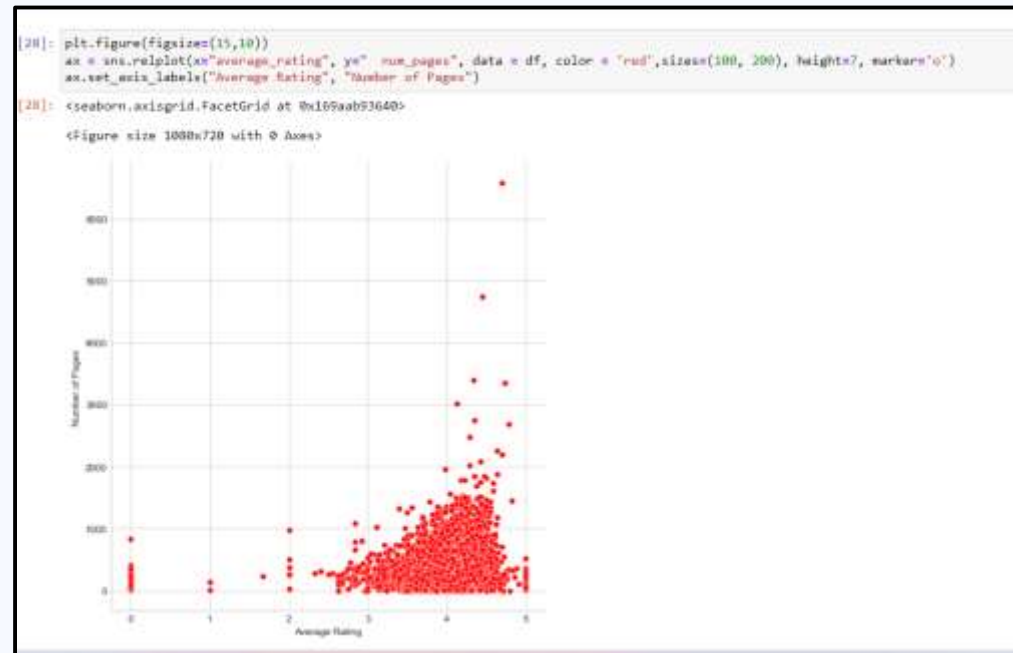# DATA VISUALIZATION

# DATA VISUALIZATION

# DATA VISUALIZATION

DATA
VISUALIZATION

# THE MODEL

11123 rows × 34 columns

```
2]:   #Scaling the values
      from sklearn.preprocessing import MinMaxScaler
      min_max_scaler = MinMaxScaler()
      features = min_max_scaler.fit_transform(features)

      #Building the model
      model = neighbors.NearestNeighbors(n_neighbors=6, algorithm='ball_tree')
      model.fit(features)
      dist, idlist = model.kneighbors(features)

2]:   def BookRecommender(book_name):
```

# EXAMPLE OUTPUT

```
Recommender(book_name):
itialize an empty list to store recommended book names
_list_name = []

nd the index of the input book_name in the DataFrame
_id = df2[df2['title'] == book_name].index

tract the first (and only) element from the index array
_id = book_id[0]

erate through the indices of the k-nearest neighbors for the input book_
newid in idlist[book_id]:
    # Append the title of each recommended book to the book_list_name
    book_list_name.append(df2.loc[newid].title)

turn the list of recommended book names
rn book_list_name

he BookRecommender function with a specific book name
s = BookRecommender('Harry Potter and the Half-Blood Prince (Harry Potte

y the list of recommended book names
s[1:]
```

```
Potter and the Order of the Phoenix (Harry Potter  #5)',
llowship of the Ring (The Lord of the Rings  #1)',
Potter and the Chamber of Secrets (Harry Potter  #2)',
Potter and the Prisoner of Azkaban (Harry Potter  #3)',
ghtning Thief (Percy Jackson and the Olympians  #1)']
```

# OUTPUT WITH BOOK COVER