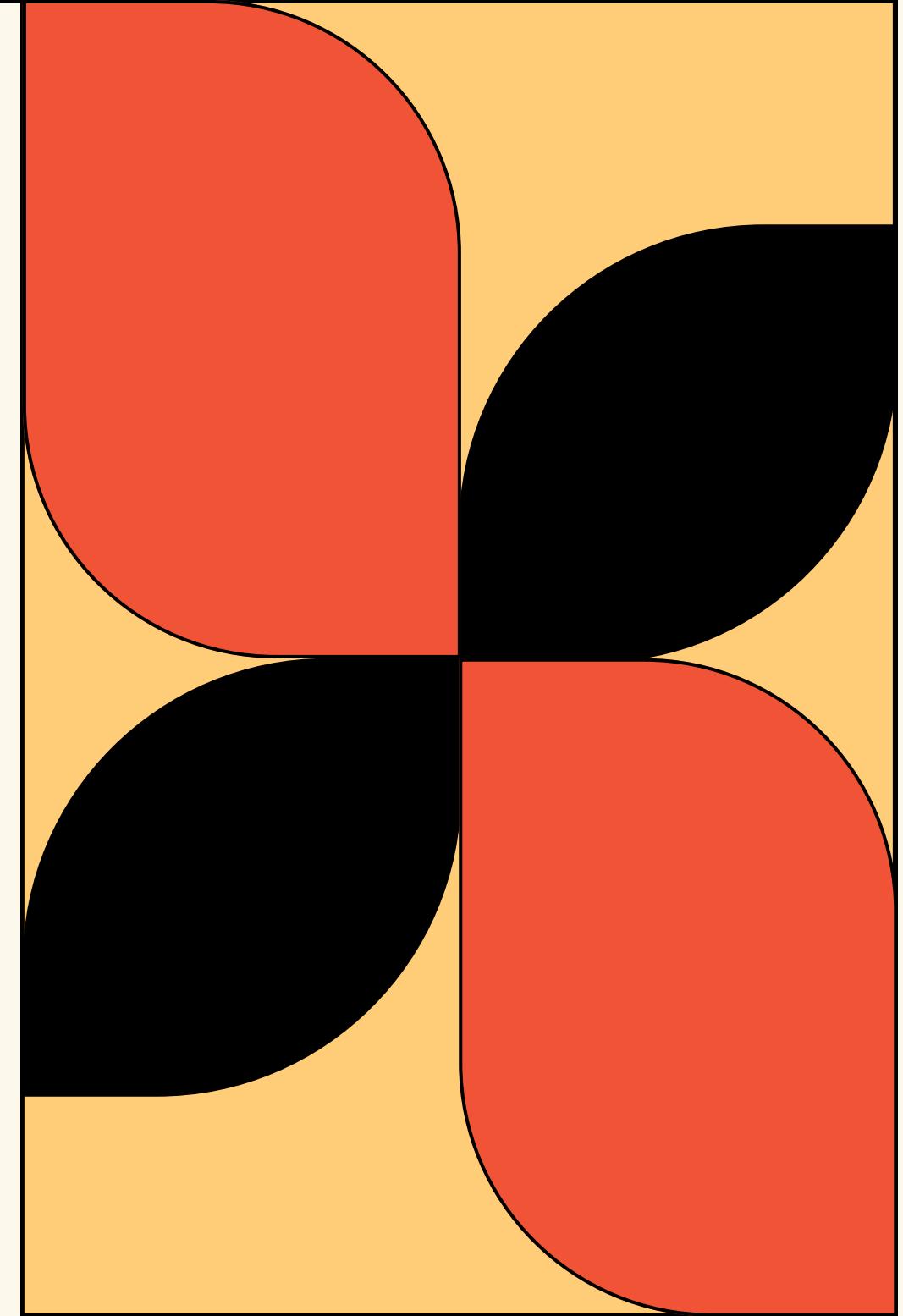




Can You Judge a Book by Its Cover?

Final Project by: Serena Lee and Sunny Yu

DATASCI 112 - WINTER 2023



Agenda

- 1 Project Goals
- 2 Data Collection & Cleaning
- 3 Data Exploration
- 4 Machine Learning



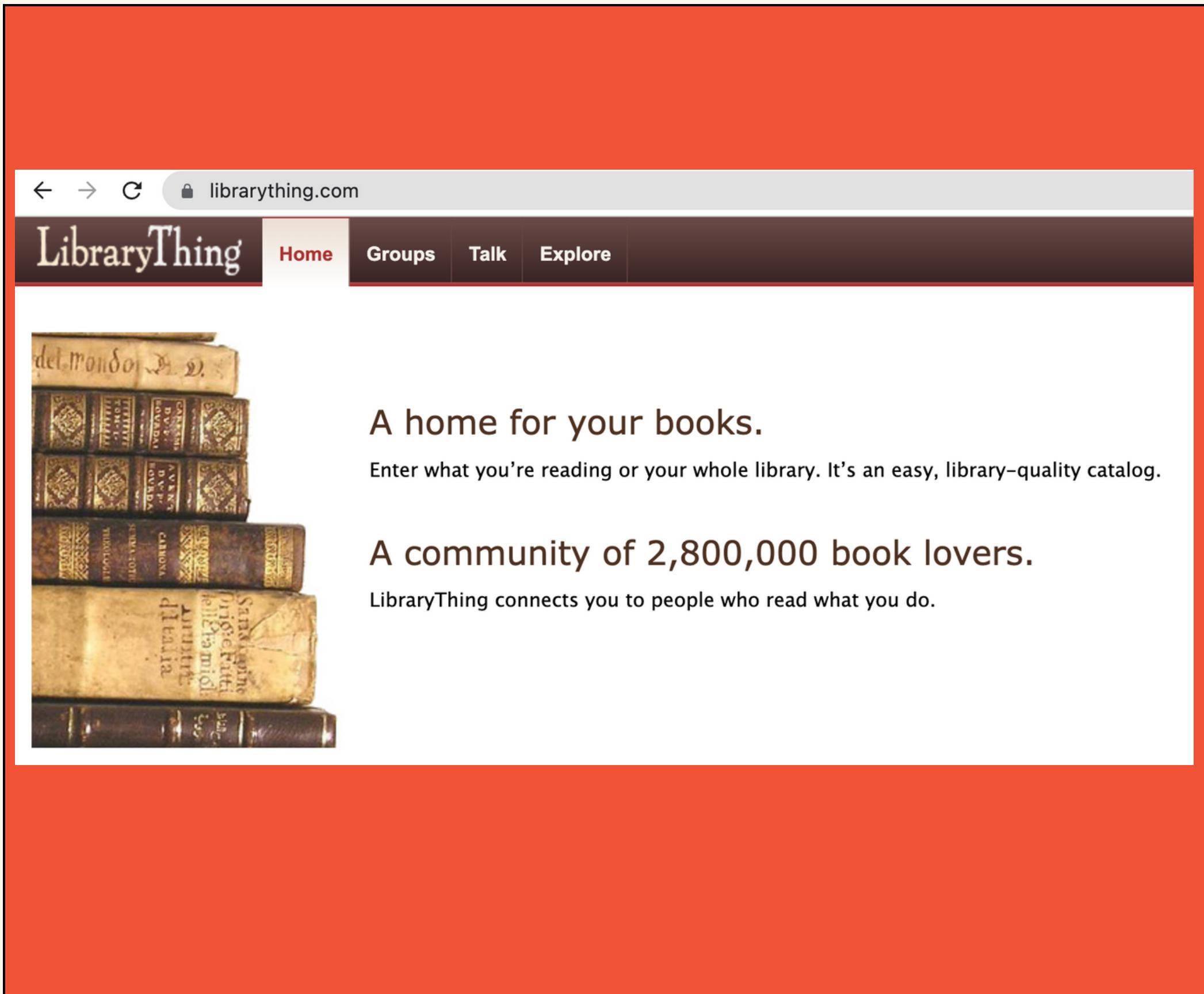
Aim 1: train a regression model to predict a book's rating based on surface-level features

- Book cover
- Blurb
- First words
- Genre
- etc...



Aim 2:
train a clustering
model to offer book
recommendations

DATA COLLECTION



- Primary source: LibraryThing
- Web scraping (x4)
 - Author names (Wiki, Guardian, LibraryThing)
 - Book page URLs
 - Book metadata
 - Book covers
- Merged multiple dataframes

DATA COLLECTION

The screenshot shows a LibraryThing book page for "Harry Potter and the Sorcerer's Stone" by J.K. Rowling. The page includes the book cover, author information, member statistics, a brief description, recently added users, and a list of tags.

Members: 123,357

Reviews: 1971

Popularity: 1

Average rating: ★★★★½ (4.29)

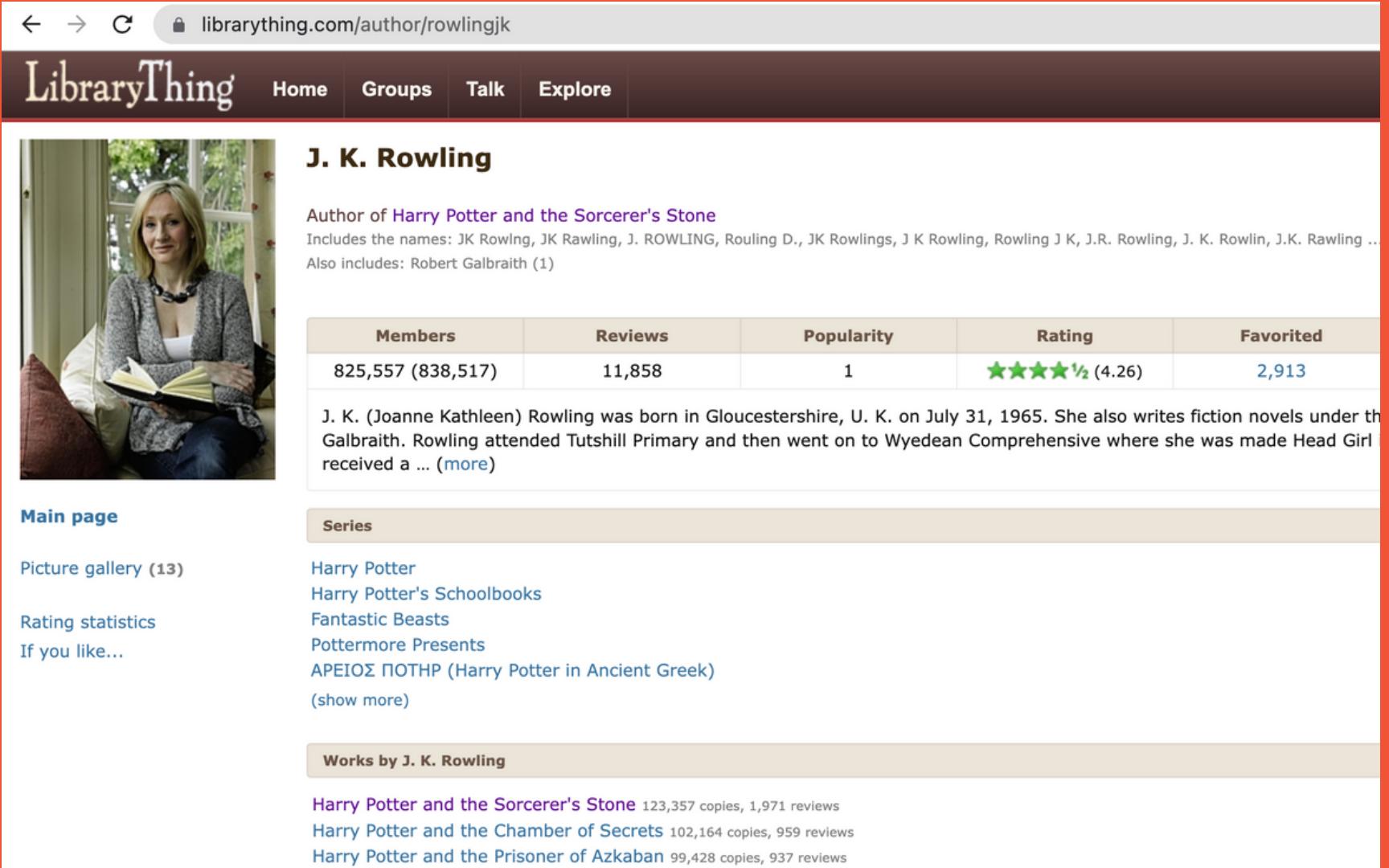
Description: Rescued from the outrageous neglect of his aunt and uncle, a young boy with a great destiny proves his worth while attending Hogwarts School of Witchcraft and Wizardry.

Recently added by: abby28472937, KristenBrinksman, chengcrickmore, jessrcarpinelli, ChristinaRobertson, Britishrose92, JeradBumblebeefromhell

Tags: adventure, boarding school, British, children, children's fiction, children's literature, England

- Primary source: LibraryThing
- Web scraping (x4)
 - Author names (Wiki, Guardian, LibraryThing)
 - Book page URLs
 - Book metadata
 - Book covers
- Merged multiple dataframes

DATA COLLECTION



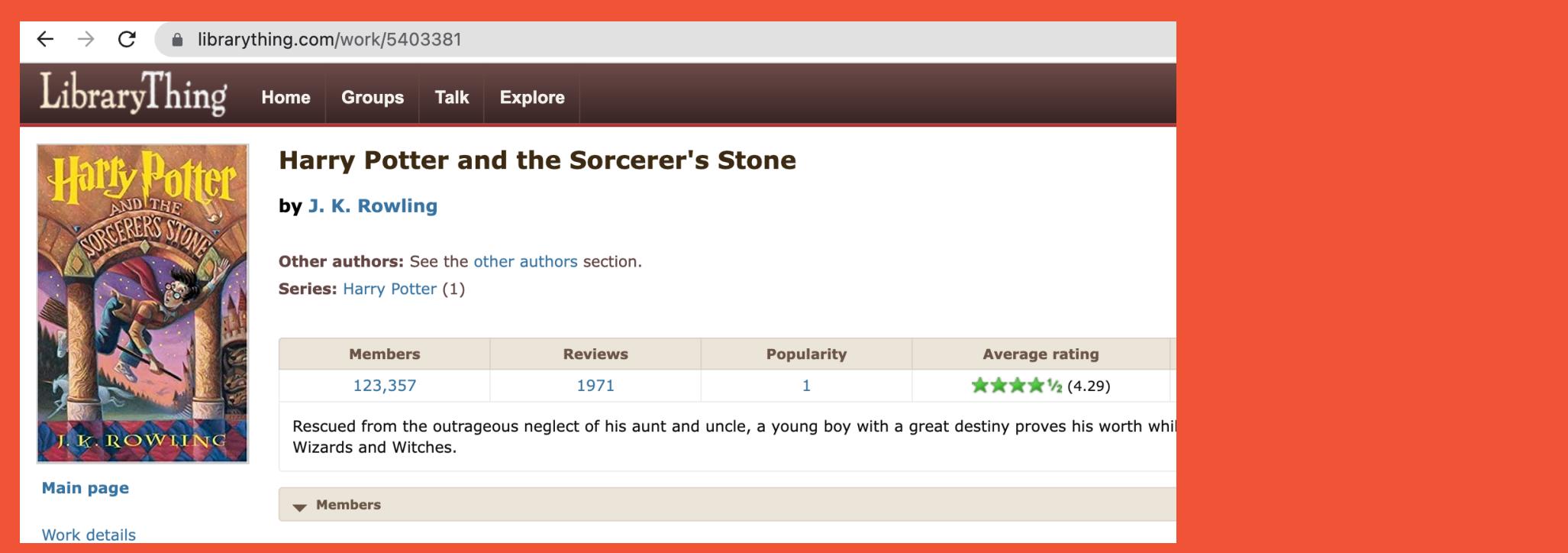
The screenshot shows the LibraryThing author profile for J.K. Rowling. At the top, there's a navigation bar with links for Home, Groups, Talk, and Explore. Below the navigation is a photo of J.K. Rowling sitting on a couch, holding an open book. The title "J. K. Rowling" is displayed in bold. A brief bio states she was born on July 31, 1965, in Gloucestershire, U.K., and writes fiction novels under the name Robert Galbraith. Below the bio is a table showing statistics: Members (825,557), Reviews (11,858), Popularity (1), Rating (★★★½ (4.26)), and Favorited (2,913). Further down, there's a section titled "Series" listing books like Harry Potter, Harry Potter's Schoolbooks, Fantastic Beasts, Pottermore Presents, and APEΙΟΣ ΠΟΤΗΡ (Harry Potter in Ancient Greek). At the bottom, a section titled "Works by J. K. Rowling" lists three books: Harry Potter and the Sorcerer's Stone, Harry Potter and the Chamber of Secrets, and Harry Potter and the Prisoner of Azkaban.

step 1: scrape 1613 author names from Wikipedia, The Guardian, and LibraryThing to get "author_names"

step 2: append author names to LibraryThing dir, loop through the author URLs to scrape their books to get "book_urls"

- Primary source: LibraryThing
- Web scraping (x4)
 - Author names (Wiki, Guardian, LibraryThing)
 - Book page URLs
 - Book metadata
 - Book covers
- Merged multiple dataframes

DATA COLLECTION



The screenshot shows a LibraryThing book page for "Harry Potter and the Sorcerer's Stone" by J. K. Rowling. The page includes the book cover, author information, member statistics (123,357 members, 1971 reviews), and a short description. Below the main content is a modal window displaying canonical and original titles, alternative titles, publication date, and people/characters.

Members	Reviews	Popularity	Average rating
123,357	1971	1	★★★★½ (4.29)

For more help see the [Common Knowledge help page](#).

Canonical title [Harry Potter and the Sorcerer's Stone](#)

Original title [Harry Potter and the Philosopher's Stone](#)

Alternative titles [Harry Potter and the Sorcerer's Stone \(US\)](#)

Original publication date [1997-07-26](#)

People/Characters [Harry James Potter](#)

step 3: loop through book_urls to turn book metadata into a data frame

- Primary source: LibraryThing
- Web scraping (x4)
 - Author names (Wiki, Guardian, LibraryThing)
 - Book page URLs
 - **Book metadata**
 - **Book covers**
- Merged multiple dataframes

title	author	book_index	book_url	avg_rating	genre	description	publication_year	first_sentence	number_of_reviews
And Then There Were None	Agatha Christie	0	https://www.librarything.com/work/7962202	4.14	[Fiction and Literature, Mystery]	Ten houseguests, trapped on an isolated island...	1939	In the corner of a first-class smoking carriag...	500
Murder on the Orient Express	Agatha Christie	1	https://www.librarything.com/work/2742	4.07	[Fiction and Literature, Mystery]	Agatha Christie's most famous murder mystery, ...	1934	It was five o'clock on a winter's morning in S...	394
The Murder of Roger Ackroyd	Agatha Christie	2	https://www.librarything.com/work/3011	4.06	[Fiction and Literature, Mystery]	Agatha Christie's most daring crime mystery -	1926	Mrs Ferrars died on the night of the 16th-17th...	291

step 4: after looping through 119,215 URLs, we collected all the book data in separate runtimes, so we concatenated multiple data frames. We also merged the data with additional info we got from Kaggle (author identities).

step 5: scraped, read, resize image data from image URLs into a 1-D array for ML



DATA EXPLORATION

popularity of genres, authors, books over time

author information

common first words of a book and its trend

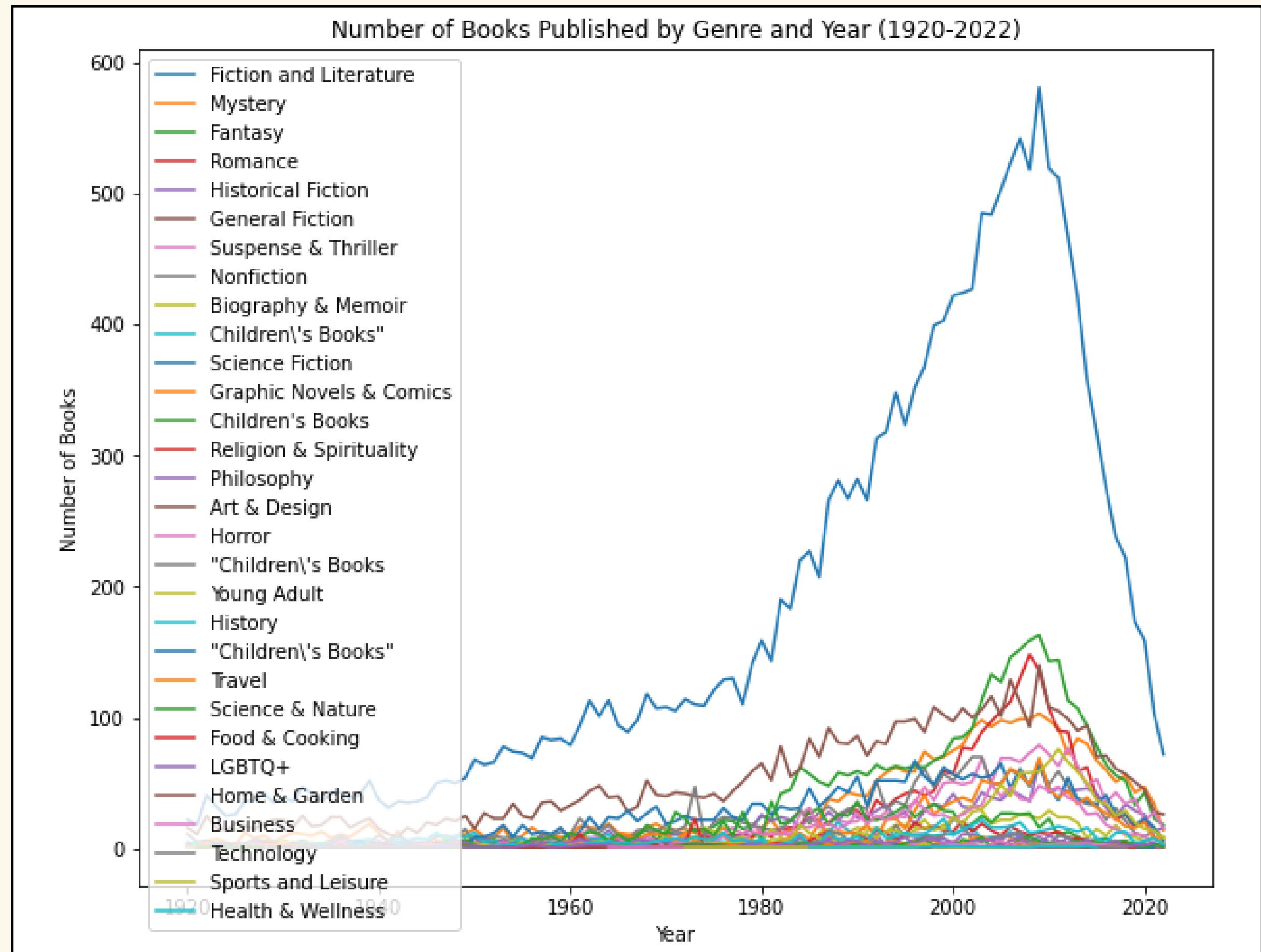
Our data comprised **24,971 books** published between **120AD and 2023** and spanned authors from **67 countries**. We wanted to explore trends across time, country, and genre.

Our exploration led us to these questions:

- 1. Where are the authors in our dataset from?** Which countries produced the most best-selling authors?
- 2. Which genres were trending at what times?**
- 3. What were the most common first words in books with ratings above 4.0+?**
- 4. What are the common first word/words that begin a book, and how do they change over time?**

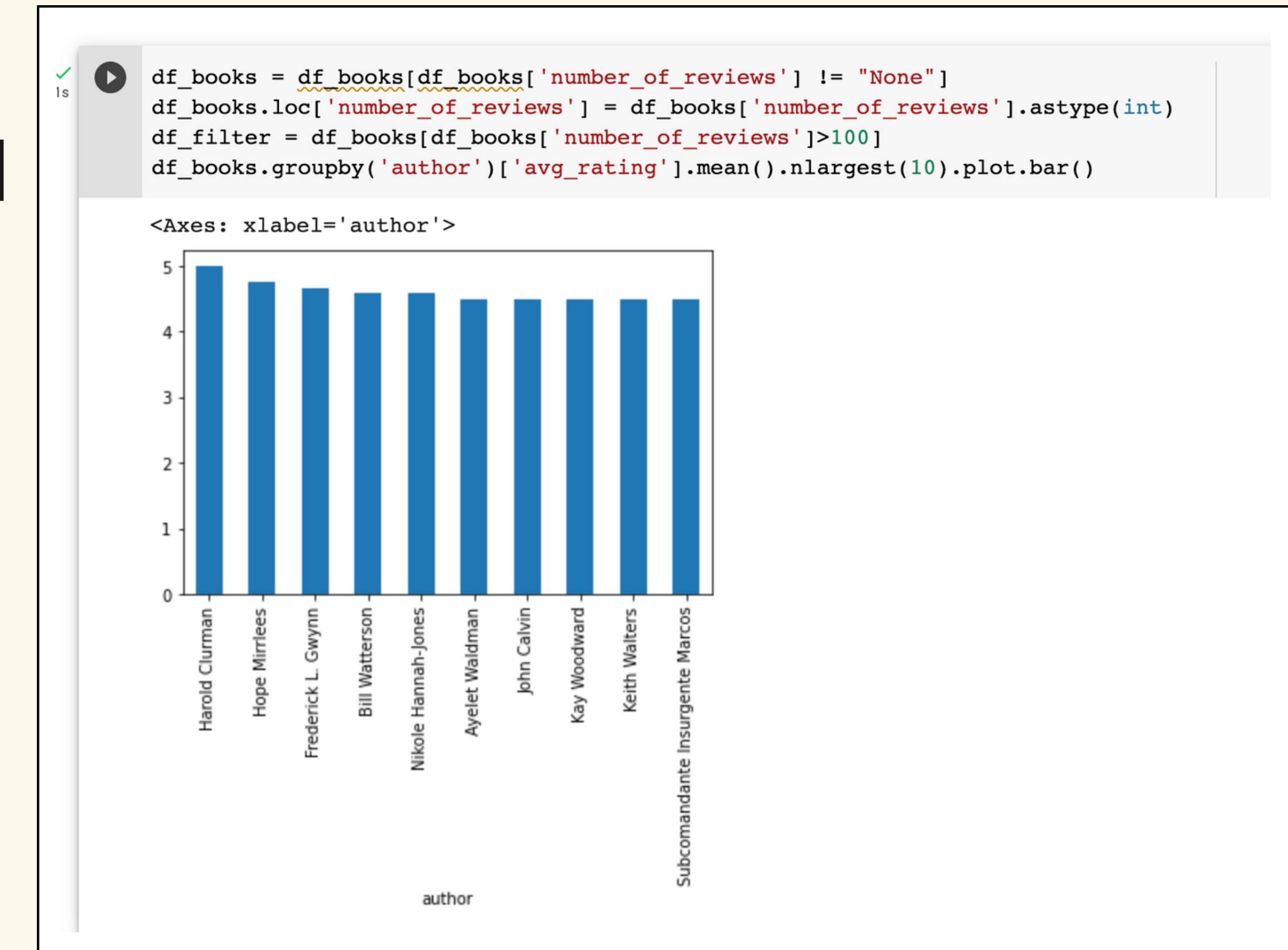
DATA EXPLORATION

POPULARITY OF GENRES



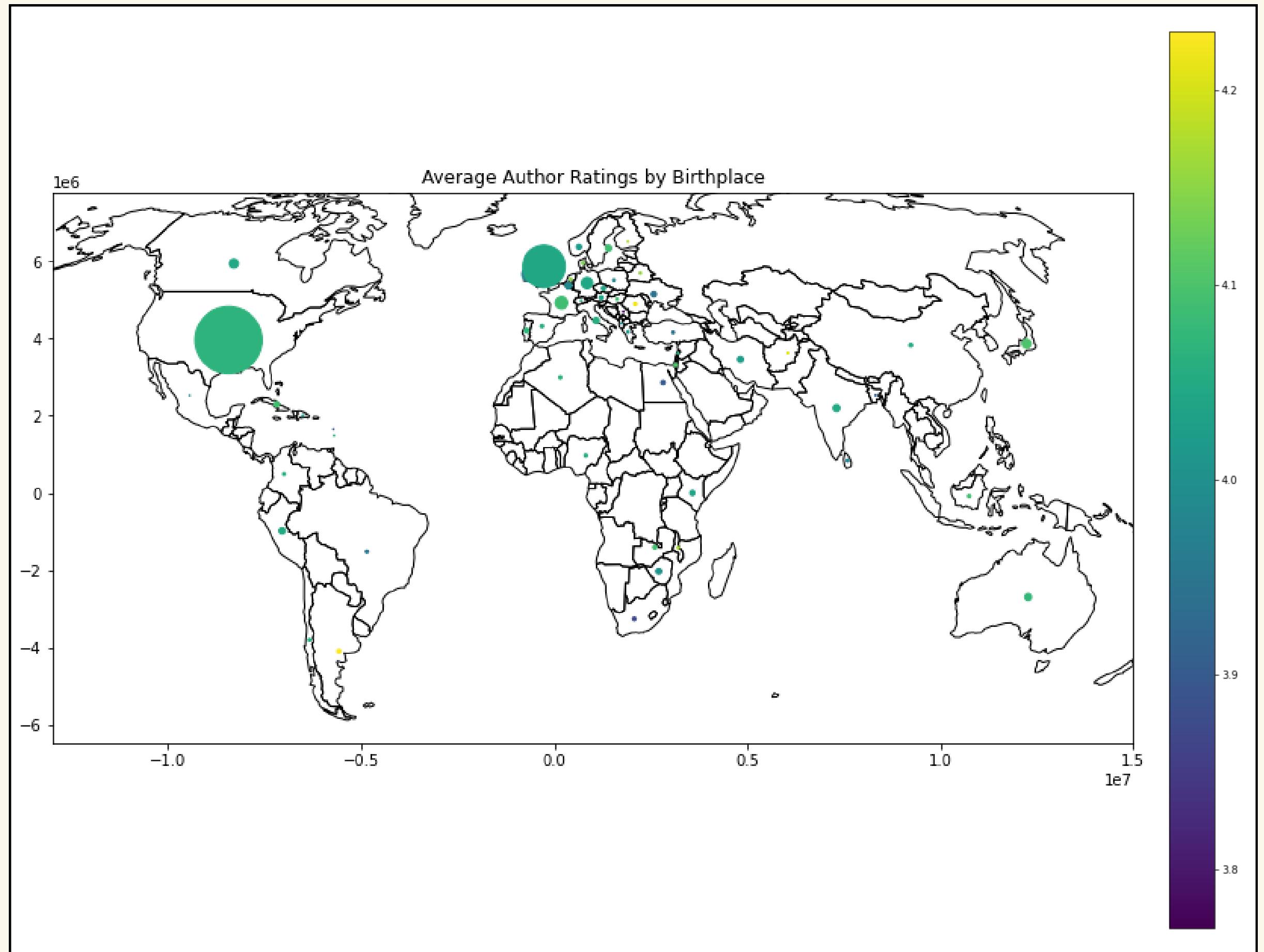
DATA EXPLORATION

POPULARITY OF AUTHORS



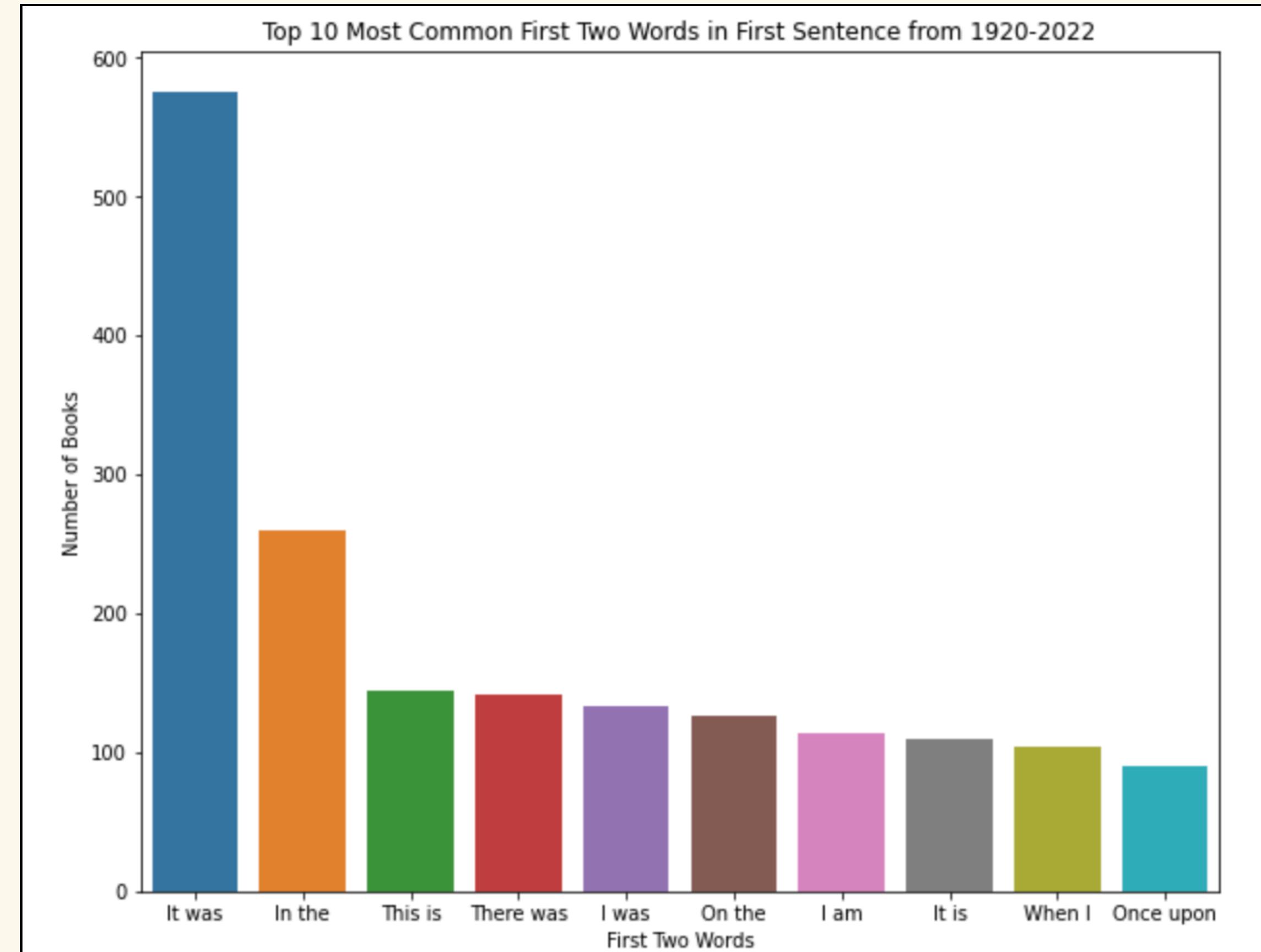
DATA EXPLORATION

AUTHOR
BIRTHPLACE



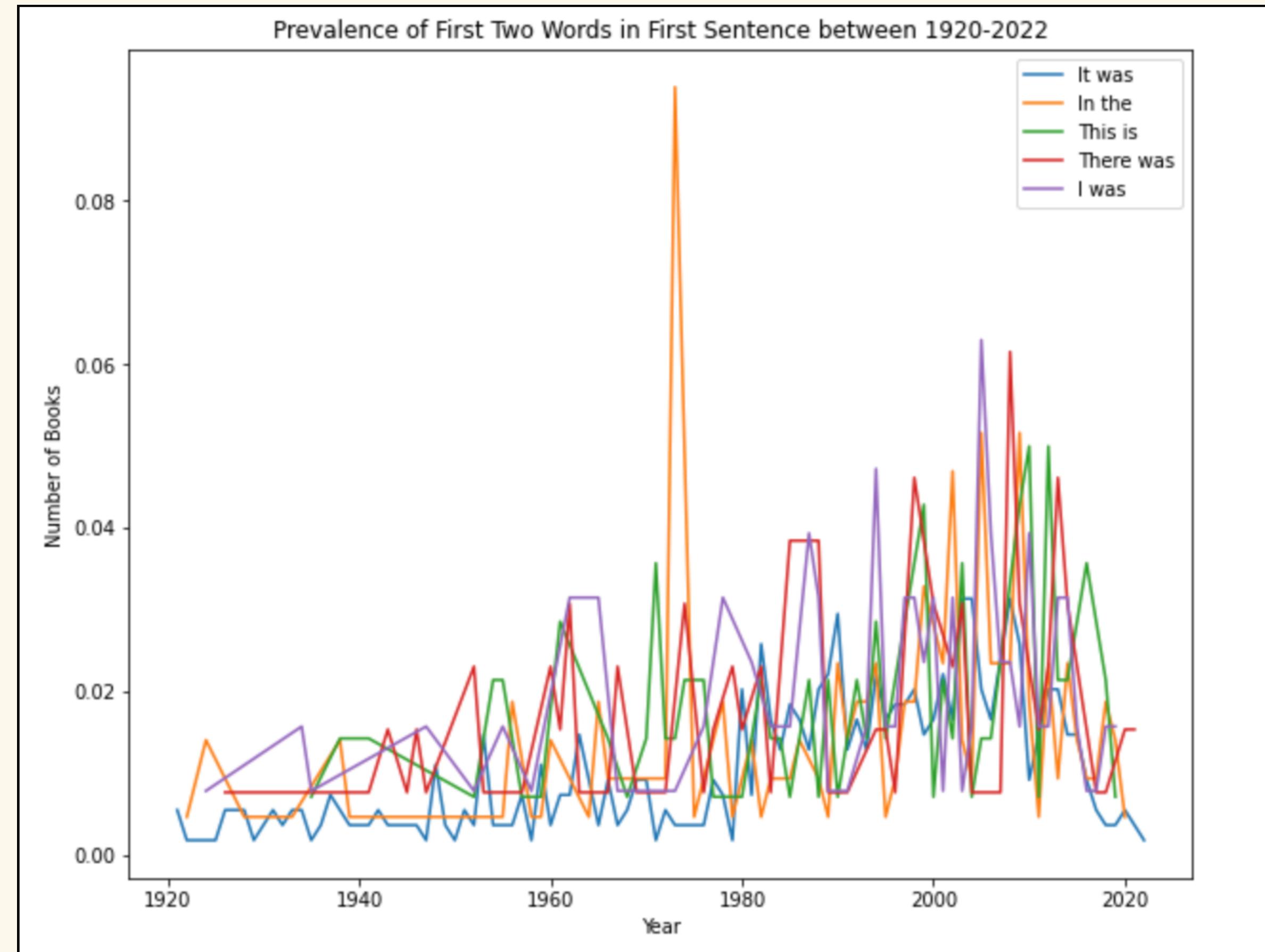
DATA EXPLORATION

*FIRST TWO WORDS:
FREQUENCY*



DATA EXPLORATION

*FIRST TWO WORDS:
TRENDS*



Machine Learning

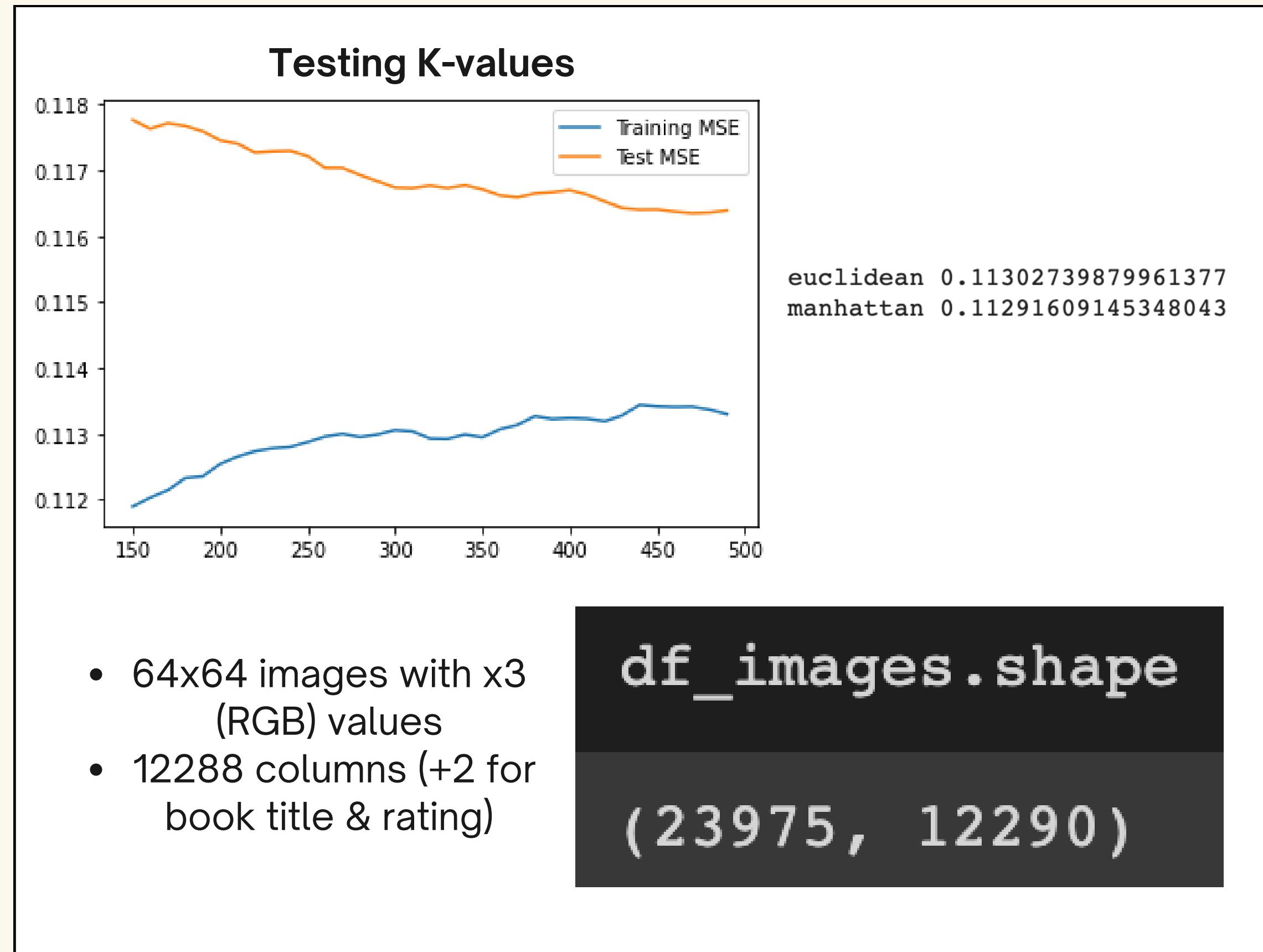
1. Predicting Book Ratings (regression)

2. Clustering similar books

Machine Learning

1. Predicting Book Ratings (regression)

- Using book covers
- KNN model of RGB pixel arrays

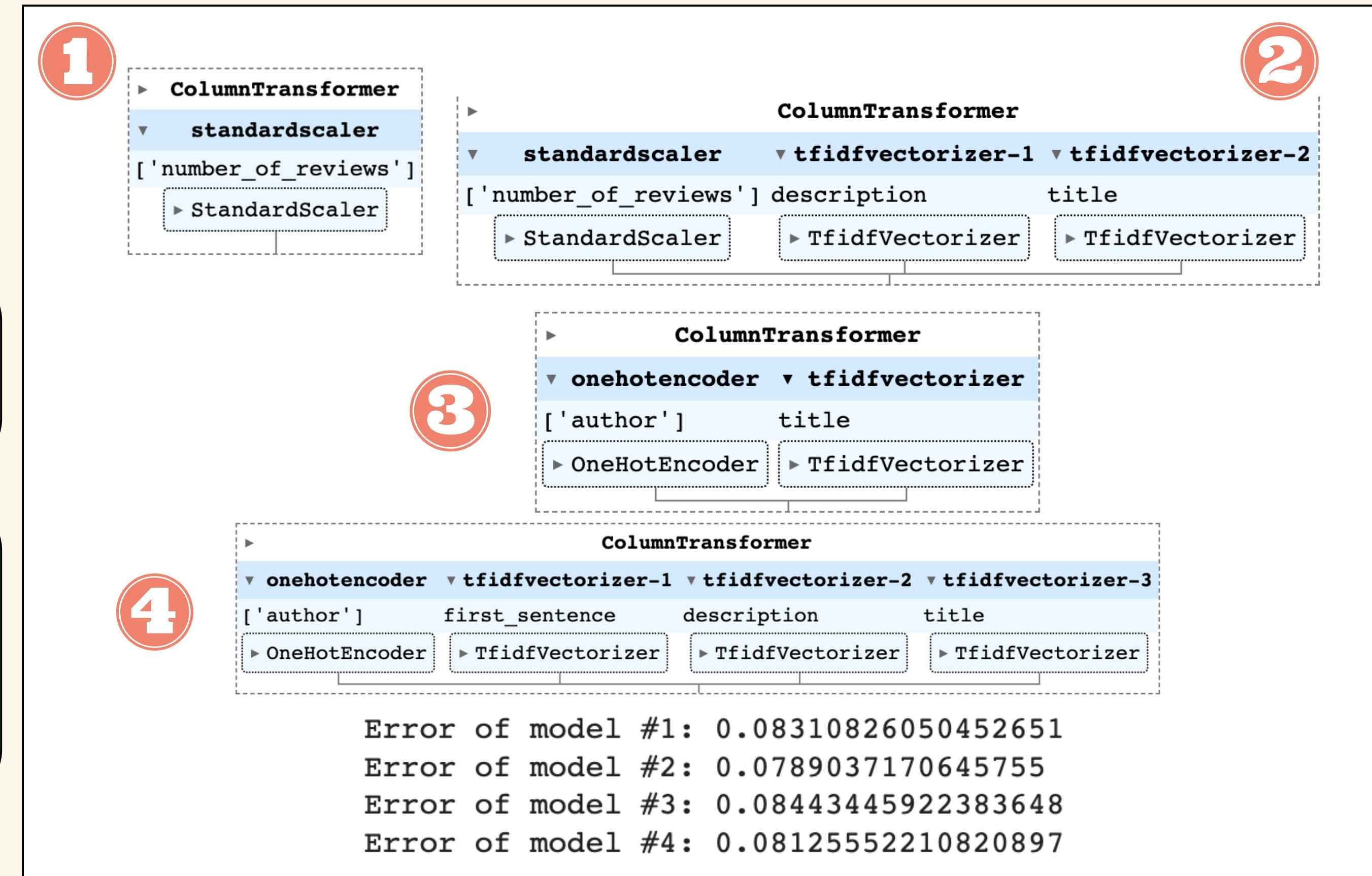


Machine Learning

COMPARISON:
Different sets of
non-image features

Best model: 0.0789

- *number of reviews*
- *description*
- *title*



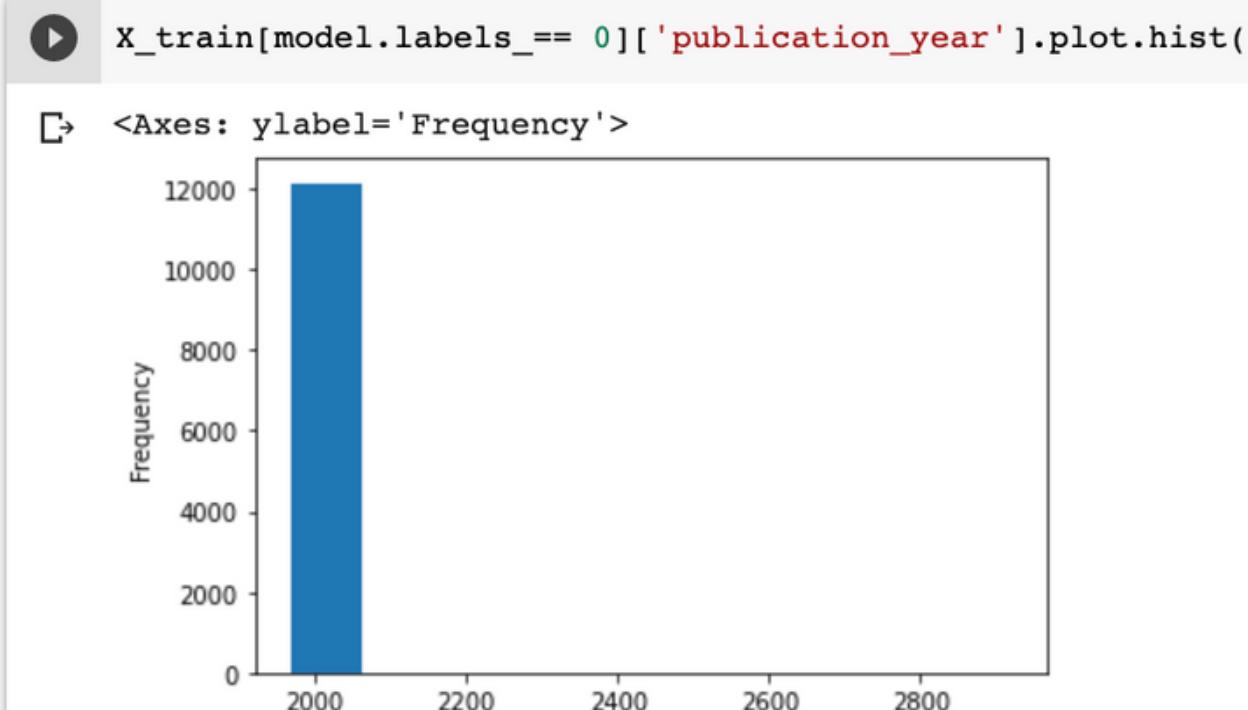
Machine Learning

2. Clustering similar books

Finding clusters

- K-means
- Hierarchical

cluster 0: modern fantasy books/fiction



```
[69] x_train[model.labels_ == 0][['author', 'genre']].value_counts()
```

author	genre	
Clive Cussler	['Fiction and Literature', 'Suspense & Thriller']	84
Stan Berenstain	["Children's Books"]	84
Margaret Weis	['Fiction and Literature', 'Fantasy']	84
M. C. Beaton	['Fiction and Literature', 'Mystery']	48
Robert B. Parker	['Fiction and Literature', 'Mystery']	47
		..
Jennifer Weiner	['Biography & Memoir']	1
	['Fiction and Literature', 'General Fiction', 'Romance', 'Mystery']	1
	['Fiction and Literature', 'Romance', 'General Fiction']	1
	['Fiction and Literature']	1
Åsne Seierstad	['Nonfiction', 'History', 'Biography & Memoir']	1
	Length: 3716, dtype: int64	

Machine Learning

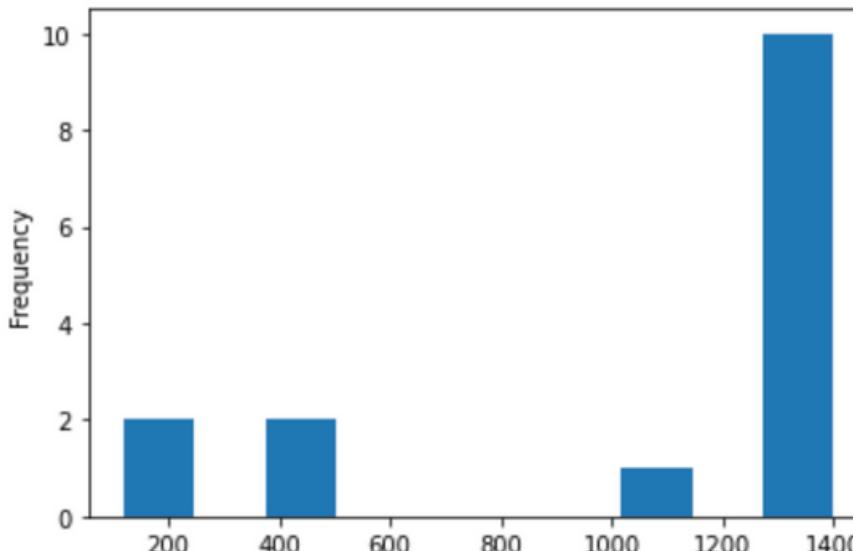
2. Clustering similar books

- Finding clusters
 - K-means
 - Hierarchical

cluster 3: ancient and early modern poets

```
▶ X_train[model.labels_ == 3]['publication_year'].plot.hist()
```

```
↳ <Axes: ylabel='Frequency'>
```



```
[75] X_train[model.labels_ == 3][['author', 'genre']].value_counts()
```

```
author           genre
Dante Alighieri  ['Fiction and Literature']
Geoffrey Chaucer  ['Fiction and Literature']
Apuleius          ['Fiction and Literature', 'General Fiction', 'Fantasy']
Geoffrey Chaucer  ['Fiction and Literature', 'Romance']
Omar Khayyam       ['Fiction and Literature']
Saint Augustine    ['Philosophy', 'Nonfiction', 'Religion & Spirituality', 'History', 'Biography & Memoir']
                     ['Philosophy', 'Nonfiction', 'Religion & Spirituality']
Suetonius          ['Nonfiction', 'History', 'Biography & Memoir']

dtype: int64
```

Machine Learning

2. Clustering similar books

- Using clusters**
- Recommend a list of similar books

Input: Any Book Name

Output: 10 Recommended Books!

```
book_recommender('Madame Bovary')
```

Customizable recommender functions

1. *similar authors*
2. *similar descriptions (TF-IDF)*
3. *similar writing style*
4. *holistic recommendation (all features)*

	title	author	book_index
2825	Willy's Pictures	Anthony Browne	2825
1927	Notes on Grief	Chimamanda Ngozi Adichie	1927
4818	Der König verneigt sich und tötet	Herta Müller	4818
9534	The Three Musketeers, Vol 1 (of 2)	Alexandre Dumas, pere	9534
4116	Madame Bovary	Gustave Flaubert	4116
2191	Invisible	Paul Auster	2191
2423	The Shut Eye	Belinda Bauer	2423
3188	The Wisdom of Father Brown	G. K. Chesterton	3188
5085	Time must have a stop	Aldous Huxley	5085
2687	Stars and Bars	William Boyd	2687



Thank you!

Book
Recommender: K-
Means Clustering

Function

```
#function for book recommendation
def get_similar_books(title):
    book_index = df_books[df_books['title'].str.strip() == title].index[0]
    cluster = clusters[book_index]
    books_in_cluster = df_books[clusters == cluster]
    return books_in_cluster.sample(10)

get_similar_books('The Great Gatsby')
```

Function Input: Any
Book Name!

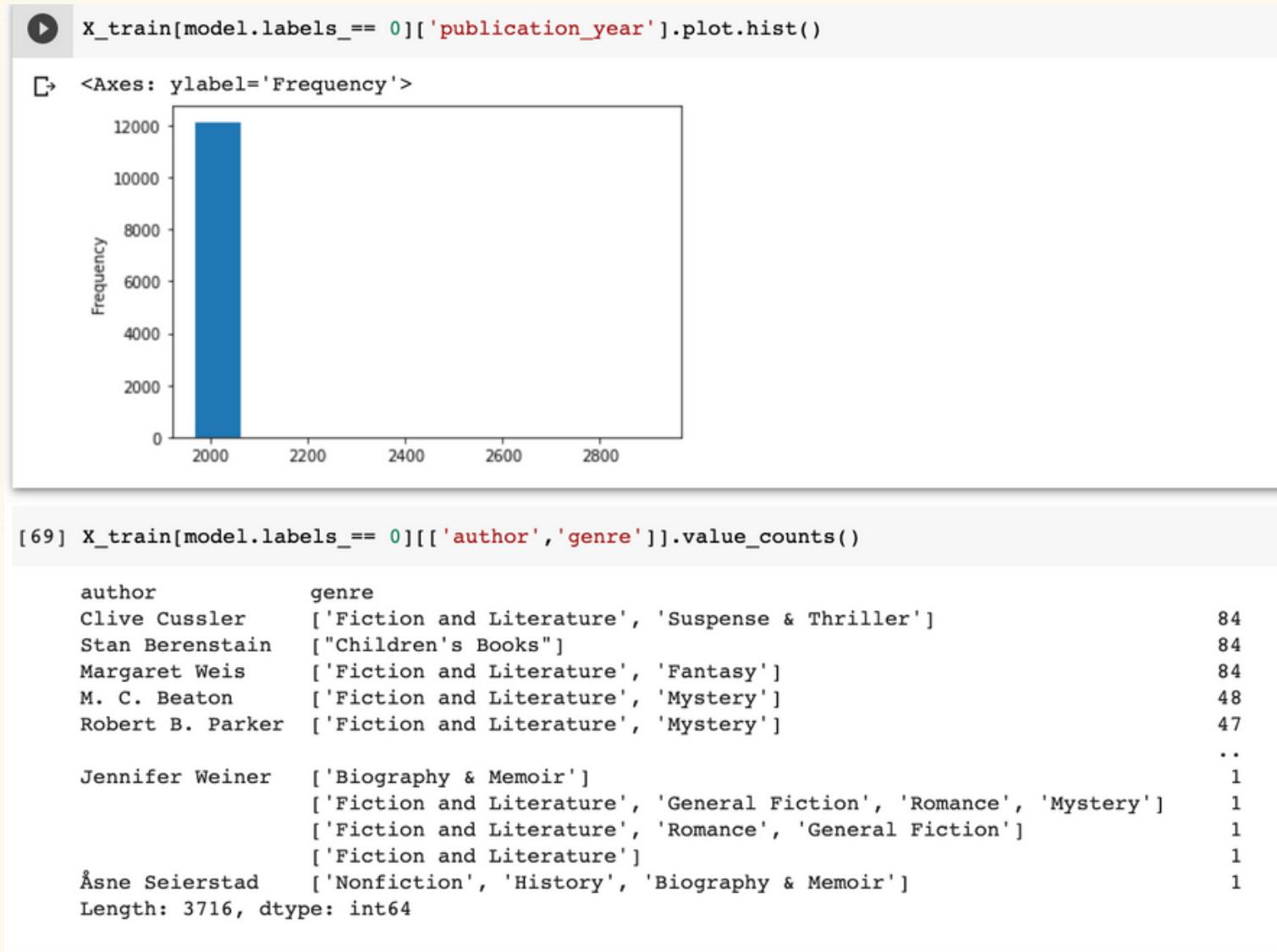
	title	author	book_
21734	Imaginary Numbers	Seanan McGuire	
12525	The Butlerian Jihad	Brian Herbert	
14399	Perfume: The Story of a Murderer	Patrick Süskind	
6319	Fall Down 7 Times Get Up 8: A Young Man's Voic...	Naoki Higashida	
24005	The Black Lyon	Jude Deveraux	
18421	Soul Harvest: The World Takes Sides	Tim LaHaye	
24601	Fever	Robin Cook	
23861	The Temptation of Forgiveness	Donna Leon	
12129	Fangirl: The Manga, Vol. 1	Rainbow Rowell	
14973	Deception Point	Dan Brown	

Function Output: A DataFrame
With 10 random
recommended books from
the same cluster

Book
Recommender: K-
Means Clustering

cluster 1: 19th century classic literature

cluster 2: mystery books



cluster 0: modern fantasy books/fiction

cluster 3: ancient Green literature



cluster 4: children's books