

## Data Collection

kaggle

Goodreads books dataset

goodreads\_books.csv

goodreads

Goodreads website scraper

Scraping of the book's reviews

Requests/BeautifulSoup

goodreads\_reviews.csv

## Data Cleaning & Restructuring

Drop column that has amazon identification number

Pandas

Add title to the "original\_title" cell if the book doesn't have one

Pandas

Fix encoding in columns that contain strings

Pandas

Drop books that don't have title, links, author, rating, number of pages, publisher, release date, genre or isbn

Pandas

Remove books with descriptions that have less than 25 characters

Pandas

Drop books that are not in the review dataset

Pandas

Used to know which books don't have reviews

Remove non-existing books from the recommended and series columns

Pandas

clean\_data.csv

Get awards

awards.csv

Get genre and votes

genre\_and\_votes.csv

Get characters

characters.csv

Drop the awards, genre\_and\_votes and characters columns (after they were used to create the new datasets)

Pandas

final\_clean\_data.csv

## Data Exploration

Analyse data

- Character Names, Descriptions and Reviews Wordclouds;
- Most present, most voted genres and description language (NLP) bar plots;
- Rating, Award Year and Publication Year distribution histograms;
- Authors and Series by number of books and Reviews text size pie charts

• Python (Pandas, Matplotlib)

Pandas

Pandas

Pandas