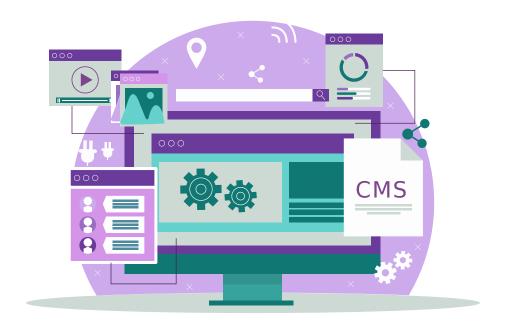


# Goodreads Books and Reviews



#### PRI, Group 2144

Inês Silva, up201806385 Mariana Truta, up201806543 Rita Peixoto, up 201806257

### **Datasets**

Books

**Books' Reviews** 



csv file



13 columns



csv file



2 columns

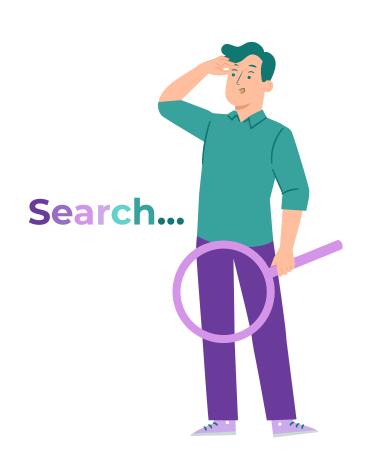


from Kaggle, originally from Goodreads

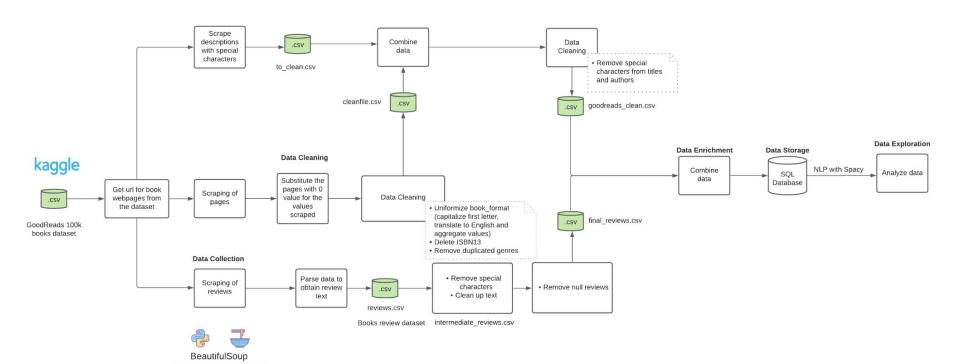
510k

lines

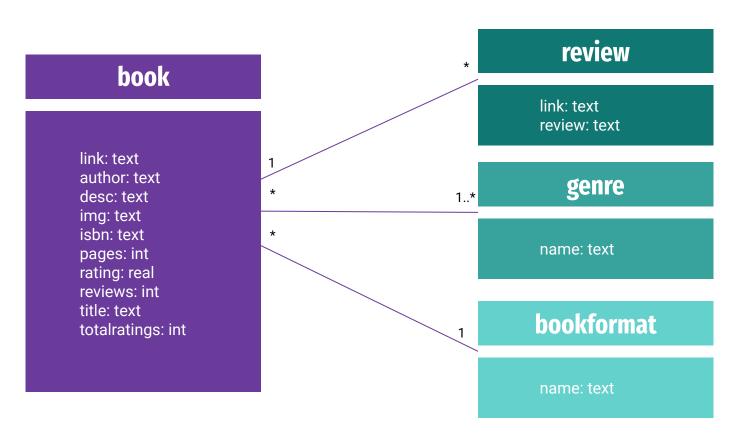
scraped from Goodreads



# **Pipeline**



# **Domain Conceptual Model**



## **Data Characterization**

**Book Format** 

Paperback

**≈ 50%** 

**Book Format** 

Hardcover

**≈ 25%** 

Rating

Average

3.83

**Review** 

Average number:

181

**Review** 

Average length (characters):

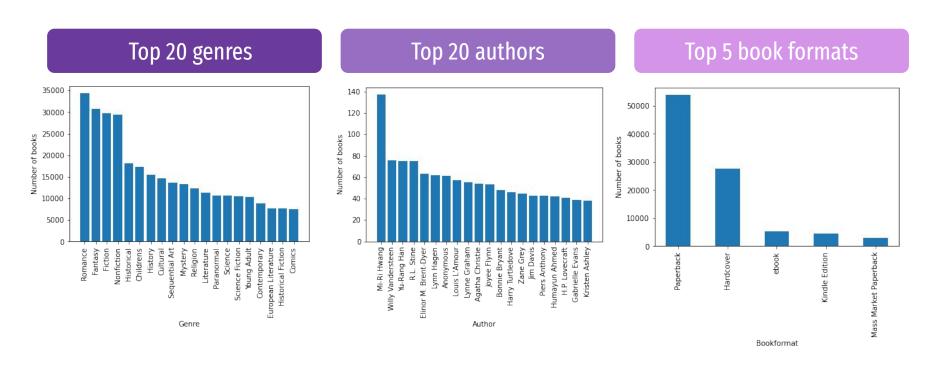
1050

**Pages** 

Average number

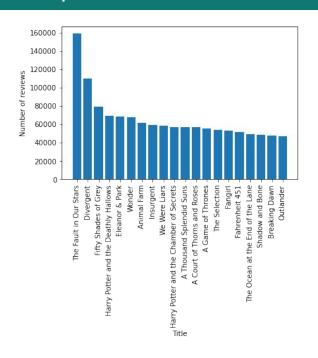
276

# **Data Characterization**

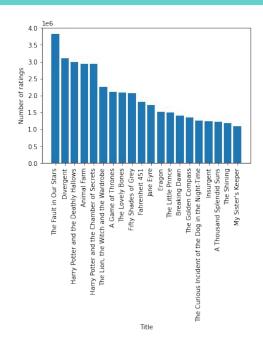


## **Data Characterization**

#### Top 20 books with most reviews



#### Top 20 books with most ratings





# **Book dataset**

author	Book format	desc	genre	img	isbn	isbn13	link	pages	rating	reviews	title	Total ratings
Charlotte Fiell,Emm anuelle Dirix	Paperback	Reveals that several hundred thousand Indians	Couture,Fa shion,Hist orical,Art, Nonfiction	https://i.g r-assets.c om/image s/S/comp ressed.ph.	190686348 2	9.78E+12	https://go odreads.c om/book/ show/100 10552-fash i	576	4.51	6	Fashion Sourceboo k 1920s	41
Week 02	Paperback	Fashion Sourceboo k - 1920s is the first book i	Politics,Hi story	https://i.g r-assets.c om/image s/S/comp ressed.ph. 	948984147	9.78E+12	https://go odreads.c om/book/ show/100 1077.Hung ar	124	4.15	2	Hungary 56	26
•••	•••	•••	•••		•••		•••		•••	•••	•••	•••



id	review
0	This collection of essays is from 1994 and the
1	When I read this book, I did not know what to
2	This was good. Very extensive. Helped in plann
3	No atendió mi expectativa, puede que se deba p

#### **Future work**



# Review Milestone 1

Revise if all the developed work is ready for the following stages



# **Building** indexes

Analyse the documents and identify their indexable components and build the indexes



## Retrieve Information

Demonstrate the indexes and retrieval process



# **Create search system**

Integration of the datasets in a search system