

Home

TV Shows

Movies

News & Popular

My list

Watchlist



DATA PREPARATION MILESTONE 1



PAULO RIBEIRO, UP201806505
PEDRO FERREIRA, UP201806506
PEDRO PONTE, UP201809694



| TABLE OF CONTENTS

01

Original Dataset

- origin
- size
- format

02

Pipeline Diagram

- Data processing flow

03

UML Diagram

- Conceptual model for the data domain

04

Data Processing

- Cleaning
- Scraping

05

Data Exploration

- Statistics Table
- Charts

06

Future Work & Improvements

- Search system
- Database storage



Home

TV Shows

Movies

News & Popular

My list

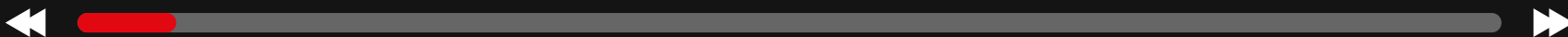
Watchlist



01

Original Dataset

Netflix dataset characterisation



Home

TV Shows

Movies

News & Popular

My list

Watchlist



02

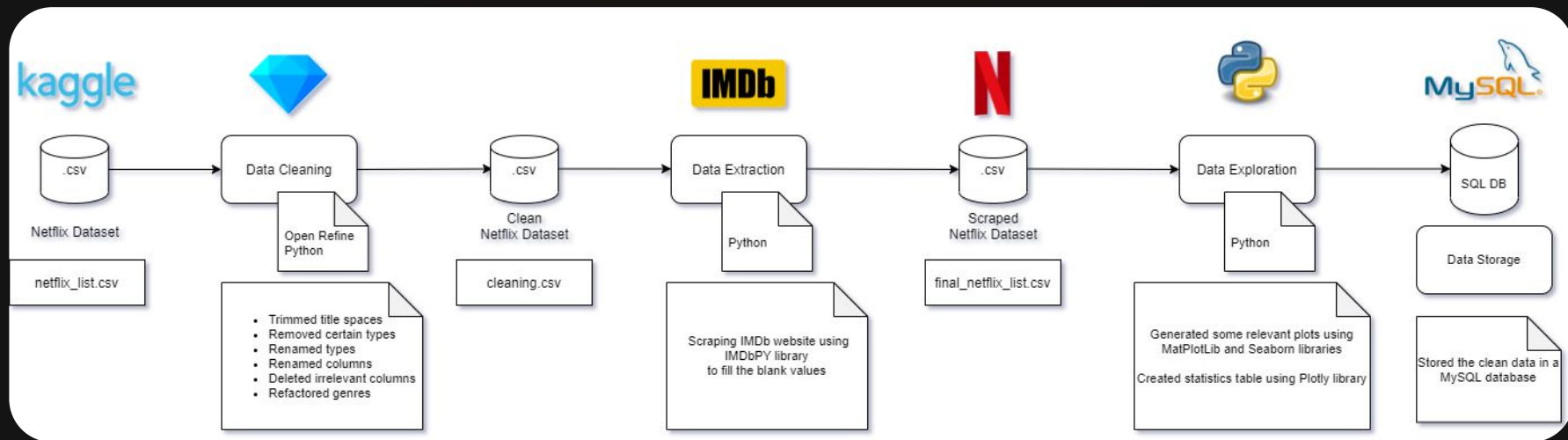
Pipeline Diagram

Data Processing flow





PIPELINE DIAGRAM



Home

TV Shows

Movies

News & Popular

My list

Watchlist



03

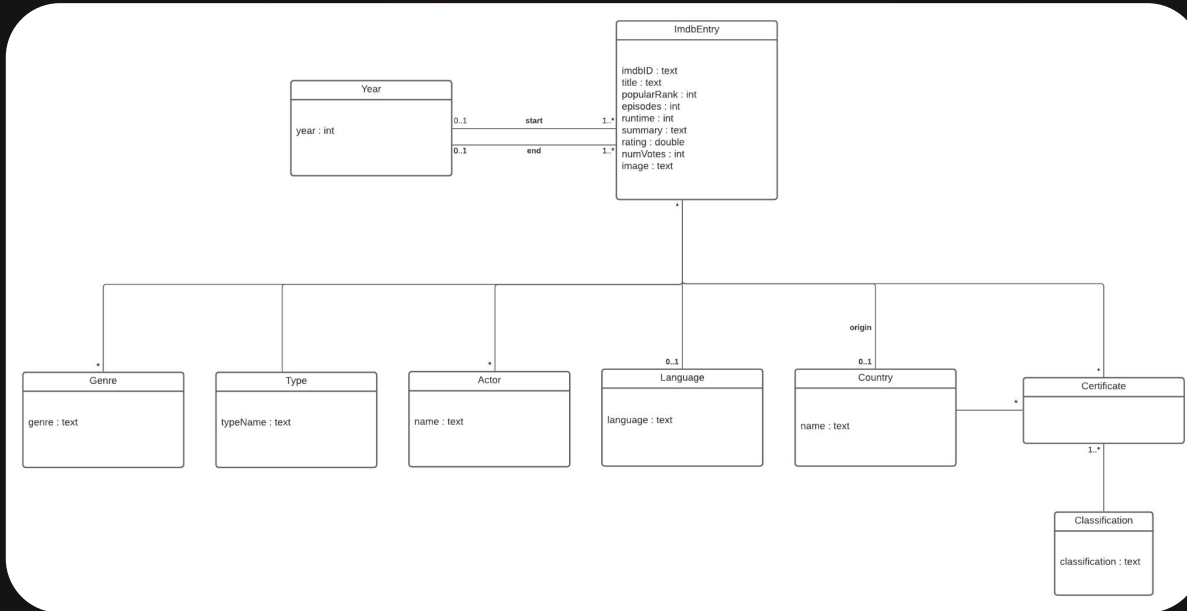
UML Diagram

Conceptual model for the data domain





| UML Diagram



Home

TV Shows

Movies

News & Popular

My list

Watchlist



04 **Data** Processing





I DATA PROCESSING



CLEANING

- Trimmed spaces
- Removed types
- Renamed types
- Renamed columns
- Deleted columns
- Refactored genres



EXTRACTION

- IMDb Scrapping using IMDbPY library



Home

TV Shows

Movies

News & Popular

My list

Watchlist



05

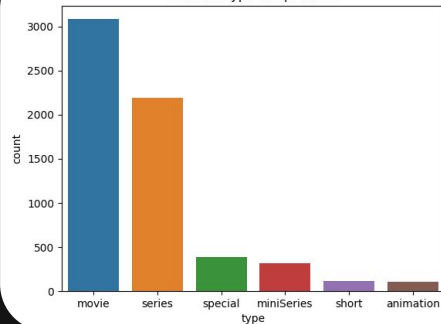
Data Exploration



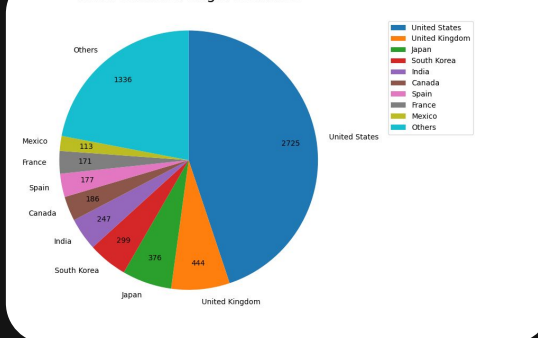
[Home](#)[TV Shows](#)[Movies](#)[News & Popular](#)[My list](#)[Watchlist](#)

| Data Exploration

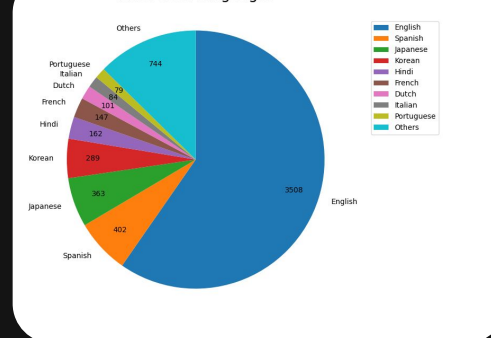
Column 'type' frequencies



Most common origin countries



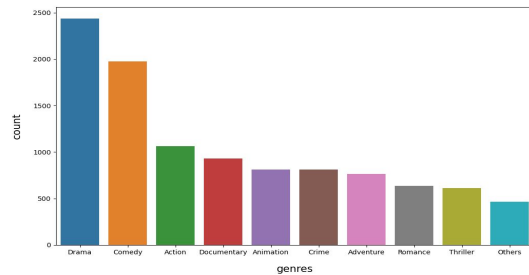
Most used languages



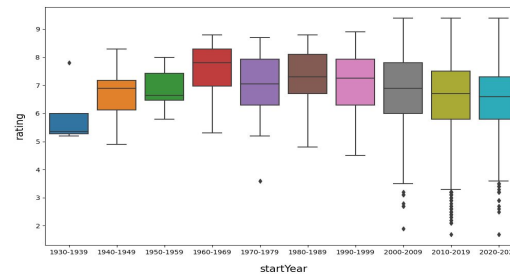
[Home](#)[TV Shows](#)[Movies](#)[News & Popular](#)[My list](#)[Watchlist](#)

| Data Exploration 2

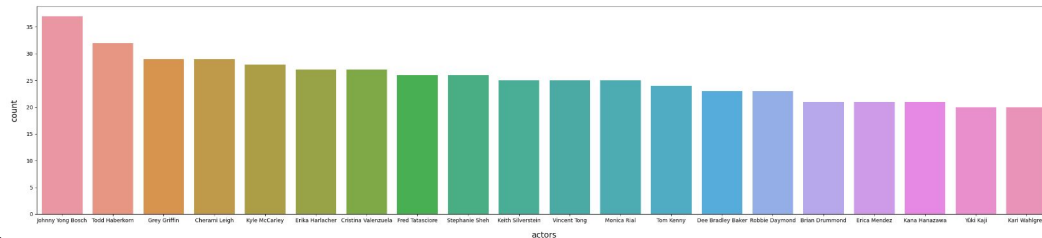
Most common genres



Evolution of rating over the decades



Actors with the most participations



[Home](#)[TV Shows](#)[Movies](#)[News & Popular](#)[My list](#)[Watchlist](#)

| Data Exploration 3

Attribute	Value
Min Rating	1.7
Max Rating	9.4
Mean Rating	6.598006524102936
Mean MiniSeries Runtime	146.73478260869564
Mean Movies Runtime	99.23572744014733
Mean Series Runtime	45.62475181998676
Mean Shorts Runtime	25.2183908045977
Mean Special Runtime	65.68115942028986
Total Genres	27
Total Languages	70
Total Origin Countries	82
Total Actors	56207



Home

TV Shows

Movies

News & Popular

My list

Watchlist



06

Future Work





I FUTURE WORK & IMPROVEMENTS

DATA STORAGE

Storing the data in a MySQL database

INFORMATION RETRIEVAL

Use of an information retrieval tool to do exploration with free-text queries

SEARCH

Use of features and techniques with the goal of improving the quality of the search results

FINAL RESULT

Our goal is an elaborate information search and retrieval system





THANKS!

Do you have any questions?



CREDITS: This presentation template was created by Slidesgo, including icons by Flaticon and infographics & images by Freepik

