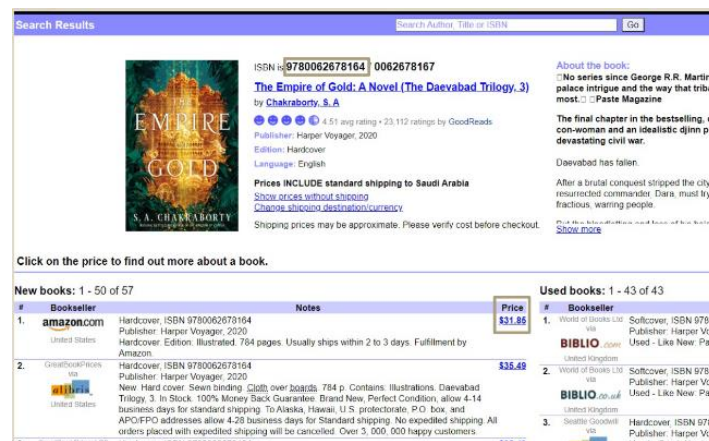**Introduction**

This is the second project for the T5 Data Science Bootcamp, which is about building a linear regression model that address a useful prediction using Scikit-learn, and using data scraped from website with Requests, BeautifulSoup, and Selenium. This project helps ReadMe firm which is an online web site that is selling books. In order to determine the price of its products they ask for our help as data scientists to study the books features to predict a suitable price.

**Design And Data**

we aim to produce a model that is able to investigate book features such as title, authors, average_rating, isbn, isbn13, language_code, num_pages, ratings_count, text_reviews_count, publication_date, and publisher. In order to determine a suitable price.

We collect the data from two deferent data sources, first source is a CSV dataset publicly available at: https://www.kaggle.com/jealousleopard/goodreadsbooks it has originally more than 11,000 rows and deferent 12 columns. We select 1298 rows and 12 columns.

The second source is the price of these books we gain it from https://www.bookfinder.com we use web scrap selenium to fetch this information. Taking advantage of selenium feature that allow the browser to interact with the web site specially in our case we make selenium searching for book prices by its isbn. Then we save the results into csv file and merge it with our previous dataset.
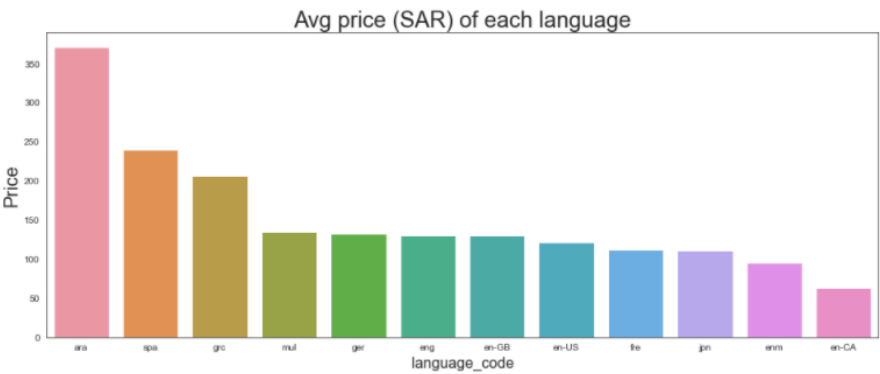


**Algorithms**

- We build a linear regression in python, and techniques such as regularization and polynomial features, adding interaction terms and dummy variables have been used.
- Rigorous model selection and evaluation has been used to select model between Linear Regression, Ridge Regression, Polynomial Regression, and Lasso Regression

**Tools**

- selenium in python (Scraping the web)
- Pandas and Numpy (Exploring the data)
- Matplotlib and Seaborn (Visualizing the data)
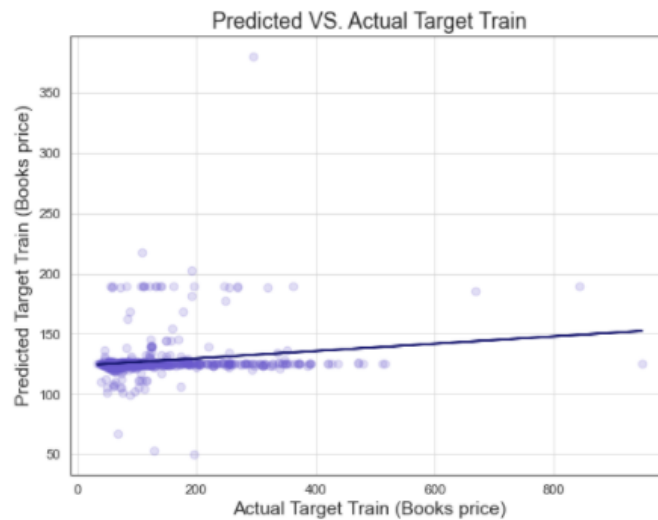- Sci-kit Learn (linear Regression model and other models)

**Communication**



- Avg price (SAR) of each language



correlation between features

Model before feature engineering



Model after feature engineering

## Cross Validation

| Regression Algorithm | Linear Regression | Ridge Regression | Polynomial Regression (degree = 2) | Lasso Regression |
|---|---|---|---|---|
| Training Score | 0.075843 | 0.044664 | 0.068357 | 0.066383 |
| Validation Score | 0.029533 | 0.00066 | 0.006979 | 0.002216 |