# Amazon Review sentiment prediction

SpringBoard Data Science Capstone Project, 2022

Hien Quang

# Objective

- Understand modern NLP techniques
- Predict if a review is negative or positive
- Set up for potential future applications such as:
    - Discover discrepancy between review sentiment and rating
    - Extract features that impactly to review's sentiment (what buyers care about most)

# Data

Source: https://nijianmo.github.io/amazon/index.html

Features:

- Reviews;
  - Ratings
  - Text
  - Helpfulness
  - Votes

- Product metadata;
  - Description
  - Category
  - Price
  - Brand
  - Image

# Data Cleaning

Huge dataset - even for the reduced version

- More than 10Gb total
- 30 json.gz files total, one for each categories

Approach:
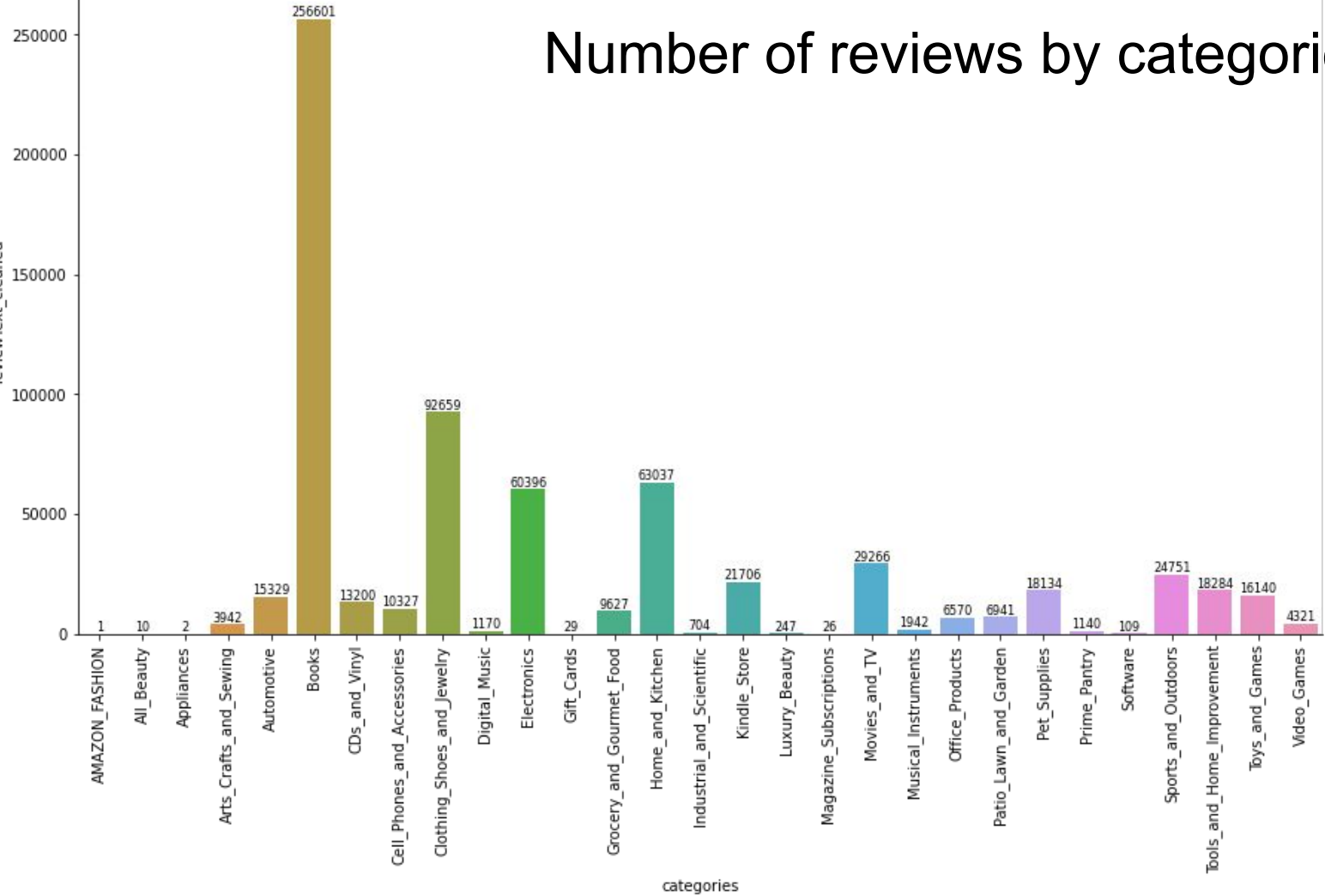
- Pyspark
- Select needed features

```
df.printSchema()
```

```
root
 |-- categories: string (nullable = true)
 |-- overall: double (nullable = true)
 |-- reviewerName: string (nullable = true)
 |-- verified: boolean (nullable = true)
 |-- vote: double (nullable = true)
 |-- reviewText_cleaned: string (nullable = true)
 |-- summary_cleaned: string (nullable = true)
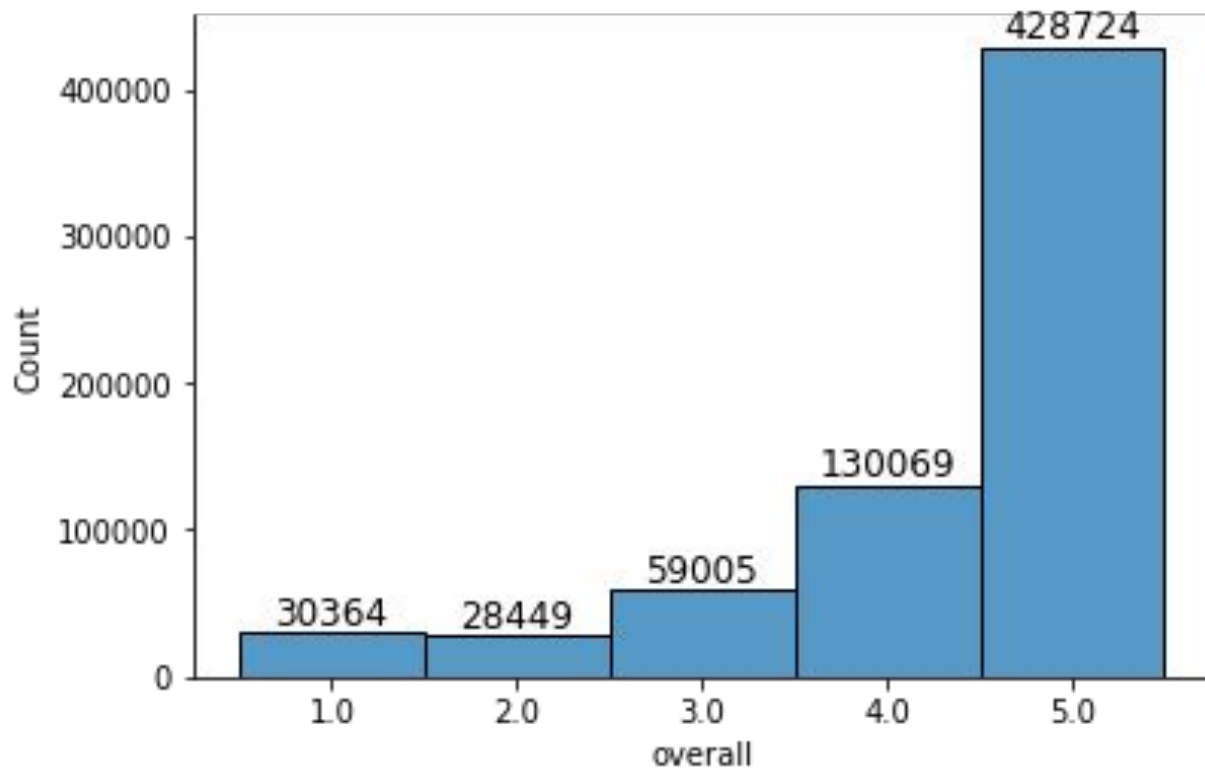```

# Exploratory Data Analysis (EDA)

Number of reviews by categories

| categories | reviewText_cleaned |
|---|---|
| AMAZON_FASHION | 1 |
| All_Beauty | 10 |
| Appliances | 2 |
| Arts_Crafts_and_Sewing | 3942 |
| Automotive | 15329 |
| Books | 256601 |
| CDs_and_Vinyl | 13200 |
| Cell_Phones_and_Accessories | 10327 |
| Clothing_Shoes_and_Jewelry | 92659 |
| Digital_Music | 1170 |
| Electronics | 60396 |
| Gift_Cards | 29 |
| Grocery_and_Gourmet_Food | 9627 |
| Home_and_Kitchen | 63037 |
| Industrial_and_Scientific | 704 |
| Kindle_Store | 21706 |
| Luxury_Beauty | 247 |
| Magazine_Subscriptions | 26 |
| Movies_and_TV | 29266 |
| Musical_Instruments | 1942 |
| Office_Products | 6570 |
| Patio_Lawn_and_Garden | 6941 |
| Pet_Supplies | 18134 |
| Prime_Pantry | 1140 |
| Software | 109 |
| Sports_and_Outdoors | 24751 |
| Tools_and_Home_Improvement | 18284 |
| Toys_and_Games | 16140 |
| Video_Games | 4321 |

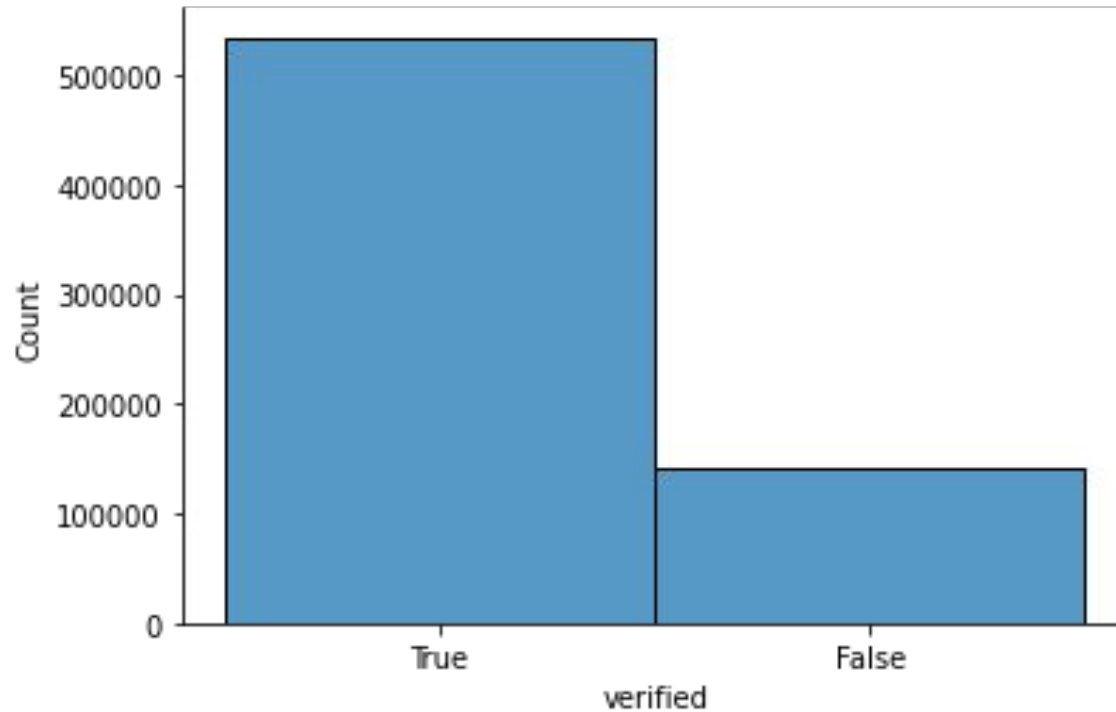# Distribution of overall score



```
count    112731.000000
mean          8.483682
std          23.656996
min           2.000000
25%           2.000000
50%           3.000000
75%           7.000000
max         999.000000
```

# Distribution of verified reviews

# Modeling

# Processed text

```python
processed_reviews = []

for review in range(0, len(X)):
    # Remove all the special characters
    processed_review = re.sub(r'\W', ' ', str(X[review]))

    # remove all single characters
    processed_review = re.sub(r'\s+[a-zA-Z]\s+', ' ', processed_review)

    # Substituting multiple spaces with single space
    processed_review = re.sub(r'\s+', ' ', processed_review, flags=re.I)

    processed_reviews.append(processed_review)
```

# TF-IDF

```
tfidfconverter = TfidfVectorizer(
    max_features=500, min_df=5, max_df=0.7,
    stop_words=stopwords.words('english'))

X = tfidfconverter.fit_transform(processed_reviews).toarray()
```

# Random Forest - result

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| negative | 0.76 | 0.31 | 0.44 | 11609 |
| positive | 0.93 | 0.99 | 0.96 | 111913 |
|  |  |  |  |  |
| accuracy |  |  | 0.93 | 123522 |
| macro avg | 0.85 | 0.65 | 0.70 | 123522 |
| weighted avg | 0.92 | 0.93 | 0.91 | 123522 |

0.9261265199721507