The background is a dark navy blue. In the top-left corner, there are two overlapping geometric shapes: a blue parallelogram and a light green parallelogram. In the bottom-left corner, there is a circular inset showing a detailed, grayscale image of a printed circuit board (PCB) with various electronic components. In the top-right corner, there is a grayscale image of a complex, multi-layered circuit board structure.

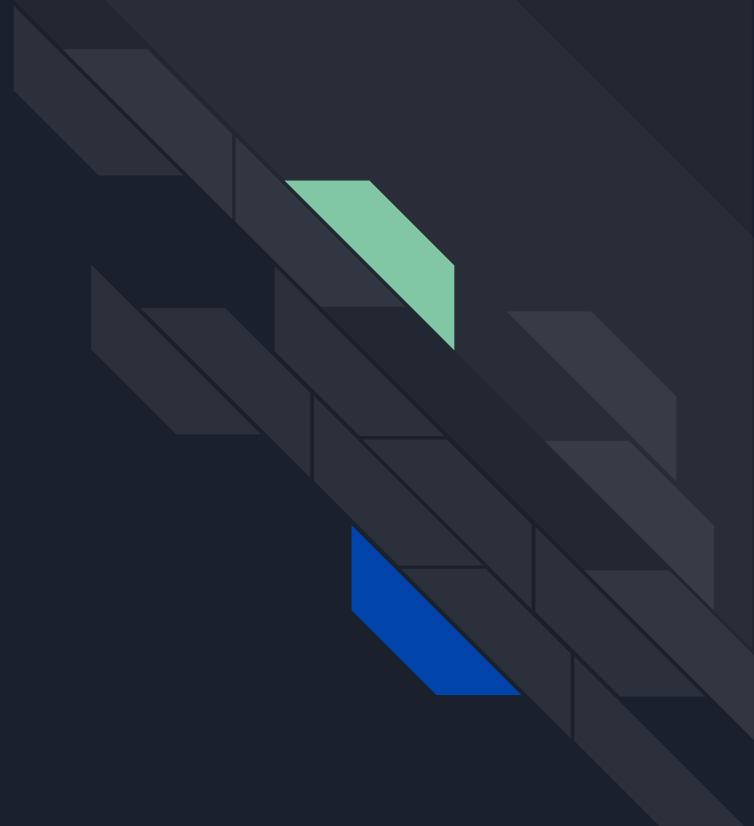
# E-commerce recommender system

Report

# The problem

What products to recommend to which customers as e-commerce retailer?

- how to ensure recommendations are 'relevant' to new customers not having history?
- and similarly, how to increase the likelihood of products with limited reviews/purchase to be seen and bought?



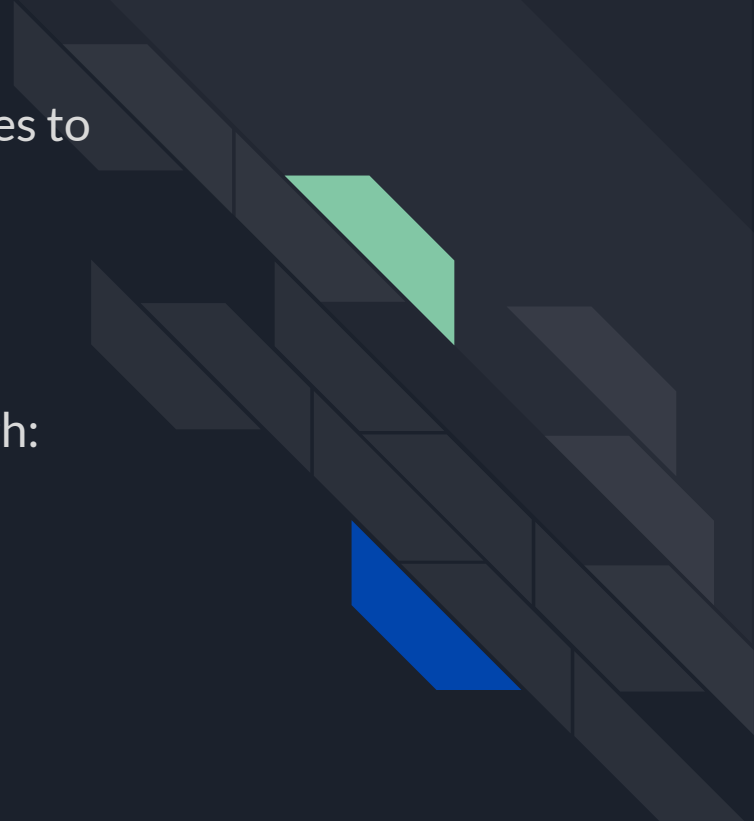
# The Approach

Develop a recommender system that learns from user/items interactions as well as from items features to recommend products users are most likely to buy

## Benefits

The recommender system will increase sales through:

- Increased basket size (value per purchase)
- Increased sales for new products
- Higher conversion rate for new customers
- Overall better experience leading to higher customers turnover





# The Data

Dataset source: Kaggle

(<https://www.kaggle.com/arashnic/marketing-bias-dataset>)

Description: the dataset contains attributes about products sold on ModCloth Amazon.

ModCloth is an indian clothing brand selling online through Amazon

Data columns includes item/user interactions (rating) and some useful data regarding products description

Data size: 99893 entries and 7 columns

Unique values: 1020 items and 44783 users

Columns:

item\_id

user\_id

Rating

size

fit

user\_attr

model\_attr

category

year

split



# Data Wrangling

## Data cleaning and tidy:

- `user_id`: There is 1 missing `user_id` value and the entire row has been deleted
- `Timestamp`: the entire column has been dropped as it won't be useful for any of our analysis
- `size`: There are about 22% missing values. We applied the most frequent value imputer strategy to populate these missing values
- `Fit`: There are about 18% of missing values. We will apply the same strategy as `size` for missing values.
- `User_attr`: there are 8% of missing values. We will apply the same strategy as `size` for missing values.
- `Brand`: There are about 74% of missing values which means this column will be useless. we will simply drop it.

	missing	cat	unique	
<code>item_id</code>	0.000000	int64	1020	
<code>user_id</code>	0.000010	object	44783	
<code>rating</code>	0.000000	int64	5	
<code>timestamp</code>	0.000000	object	14741	
<code>size</code>	0.217833	float64	9	
<code>fit</code>	0.185258	object	5	
<code>user_attr</code>	0.083760	object	2	
<code>model_attr</code>	0.000000	object	2	
<code>category</code>	0.000000	object	4	
<code>brand</code>	0.740592	object	31	
<code>year</code>	0.000000	int64	10	
<code>split</code>	0.000000	int64	3	



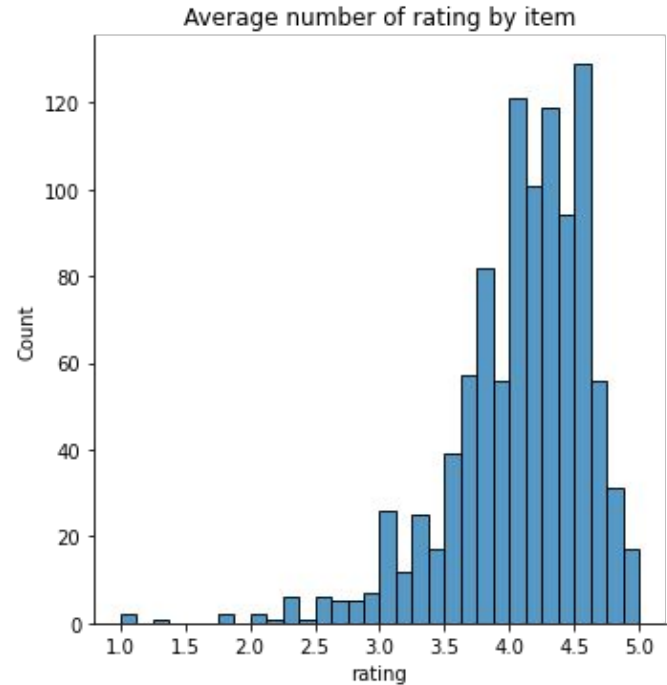
# Key Findings

- There are far more customers than products with a ratio of 43 purchases per customers
- There is considerable number of products with low ratings which is the case for the 'cold start' problem
- We have information regarding items features but not users features
- There are few 'very popular' items with high numbers of rating while others have much lower number of rating
- The data available can be used for different modelling scenario using both collaborative filtering and content base:
- Collaborative filtering: here we focus on the interactions user/item/ratings to predict products
- Contents base: we can use the items features to provide recommendations to new users or for new products



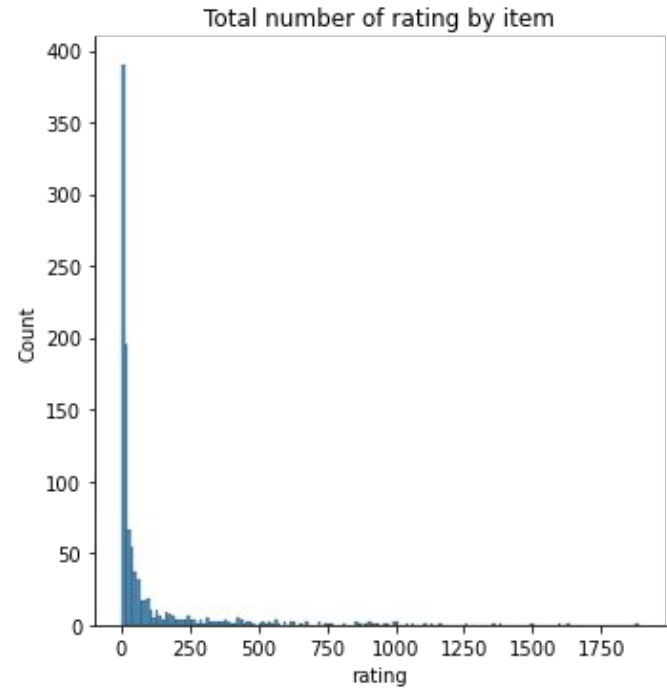
Items can be segmented into 2 categories:

- 'Popular' items' with high frequency of rating
- 'Less popular' items with few ratings



Distribution of average number of rating per item is right skewed which confirms that there are some 'very popular' items with high frequency of rating while others have much lower number of rating

## Sparse user/items interaction matrix



Most items have few number of rating which means that the item/user matrix is very sparse






# Modelling scenario

01 Collaborative Filtering using KNN and Matrix Factorization

02 Collaborative Filtering using scikit surprise library

03 Content Based Recommender

04 Hybrid Recommender




## Model selection for collaborative filtering

### Modelling using Scikit Surprise library

Algorithm	test_rmse	fit_time	test_time
<b>SVD</b>	1.056060	6.689027	0.222243
<b>KNNBasic</b>	1.080752	191.928903	11.370964
<b>NormalPredictor</b>	1.393228	0.234222	0.372735

Our final model outperformed matrix factorization with a RMSE of 1.0422 (vs 1.5 for the matrix factorization)

We used GridSearch to identify the best parameters below for our best SVD algorithm: {'n\_factors': 10, 'n\_epochs': 10, 'lr\_all': 0.005}



## Model selection for content base

Hybrid model using LightFM library

### **Modelling scenario**

LightFM was used to create a model including user/item interactions as well as items features

Loss function used is 'warp'

Number of epochs: 10

### **Evaluation**

Train auc: 98%

Test auc: 85%

Recall\_at\_100: 61%

precision\_at\_100:0%



# Recommendations

- Implement the surprise **SVD model** as go to algorithm leveraging the wealth of user/item interaction to recommend items to returning users
- Implement the **LightFM hybrid** model in parallel for items having low number of ratings as well as new users until enough data are gathered to apply the SVD model
- Finally, leverage the **content based** approach strategically for brand new products during launch phase for instance



## Further Analysis

- User features: We missed the opportunity to integrate user features in both content based and hybrid models. Information such as prior pages visited, viewed items or comments can be very useful to personalize recommendations to users.
- User profile: Initial collection of information from users can also help to instantly enrich recommendations from content based approaches. This includes information such as age, gender, height or weight.
- Real time implementation: Our implementation assumes batch feeding into the model. However, we can leverage the LightFM hybrid model to feed data partially and update it in real time. This would increase the pace of learning for our model as it could quickly adapt from user behavior



Thank you!