

Mercari Price Analysis Project

The screenshot shows the official Mercari website. At the top, there's a blue header bar with the 'MERCARI' logo, a 'Categories' dropdown, a 'Brands' dropdown, and a 'Sell on Mercari' button. To the right is a search bar with the placeholder 'What can we help you find?' and user icons for profile and cart. Below the header, a banner features the text 'Sell or buy. Almost anything.' with 'Sell now' and 'Shop brands up to 70% off' buttons. A 'Money Back Guarantee' badge is also present. To the right of the text is a collage of various items for sale, including a hat, a bag, sunglasses, and shoes. Below the banner are three circular icons: 'Sell it.' (phone icon), 'Ship it.' (truck icon), and 'Get paid.' (dollar sign icon), each with a brief description. On the far right, a 'Tell me more' link is visible.

Mercari is an e-commerce company currently operating in Japan and the United States. Their main product, called Mercari, is a marketplace app - which has grown to become Japan's largest community-powered marketplace with over JPY 10 billion in transactions carried out on the platform each month.

A problem for Mercari is price suggestions - Mercari would like to offer their sellers price suggestions based on what products they want to sell. However, this is tough, because sellers are enabled to put just about anything on the marketplace.

To help solve this problem, Mercari have put up a public dataset concerning products that have been sold on their marketplace. In this dataset, each product sold has associated properties like the name of the listing, the item description of product, the brand of the product, and more. The idea is to build an algorithm, with the dataset, that can offer price suggestions of products sellers wants to put up on their marketplace platform.

The challenge was originally posted on Kaggle:

<https://www.kaggle.com/c/mercari-price-suggestion-challenge/overview>

In this paper, the key insights of a data analysis approach to the problem is presented, and the key findings in building a linear model and applying the model to the dataset is presented.

The results indicate that the model can suggest prices for products, generally, really well in most cases, but with some difficulty suggesting prices for products that historically have been sold for very high prices

The underlying code can be found on:

<https://github.com/wildanwildan94/Mercari-Price-Analysis----Inference-Prediction-of-Products>

Mercari Price Analysis - Data Analysis - Description of dataset, item condition of products

- For an idea of the dataset, we consider the different attributes that are available and some typical associated values

- Generic Values of Each Attribute:

Name: Smashbox primer

Item Condition Id: 2

Category Name: Beauty/Makeup/Face

Brand Name: Tarte

Price: 8.0

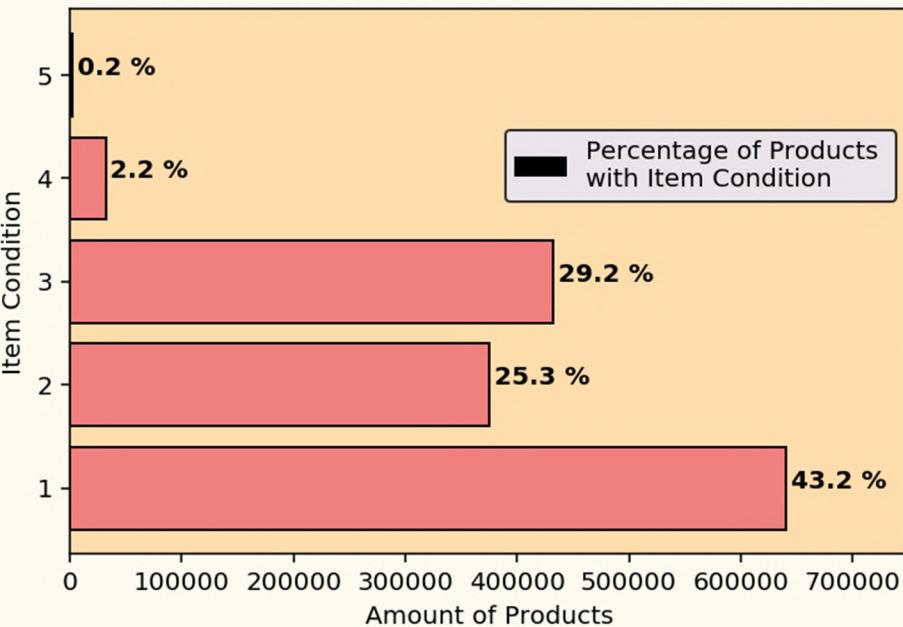
Shipping: 1

Item Description: 0.25 oz Full size is 1oz for [rm] in Sephora

- The Name is a typical, brief description of the product in question
- The Item Condition is a number representing the condition of the product in question
- The Category Name represents the category of the product
- The Brand Name is simply the brand of the underlying product, e.g. Nike
- The Price is the price the product was sold for, in the unit USD
- The Shipping is 1 if the shipping fee is paid by the seller, and 0 if it is paid by the buyer

- For each product there is an associated item condition describing the quality of the product - let us look at that

Amount of Products in each Item Condition

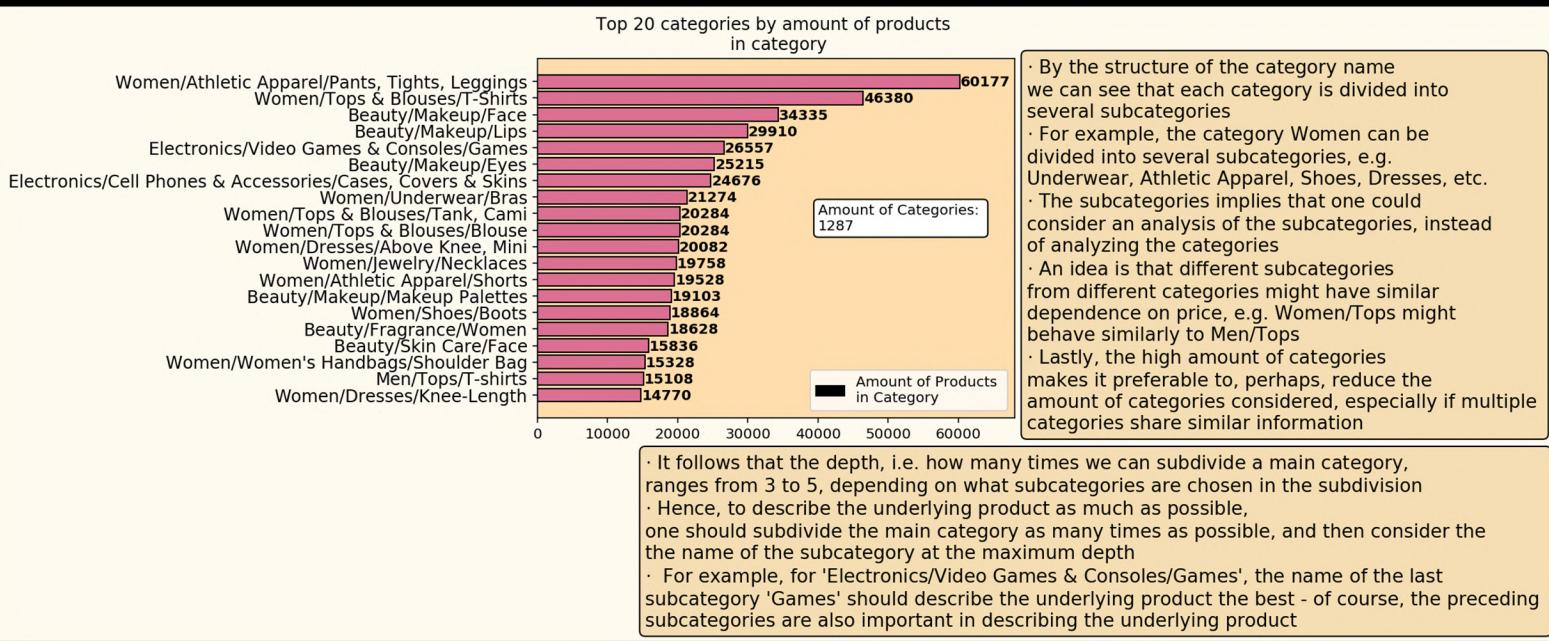


- 1: New
- 2: Almost New
- 3: Good
- 4: Fair
- 5: Poor

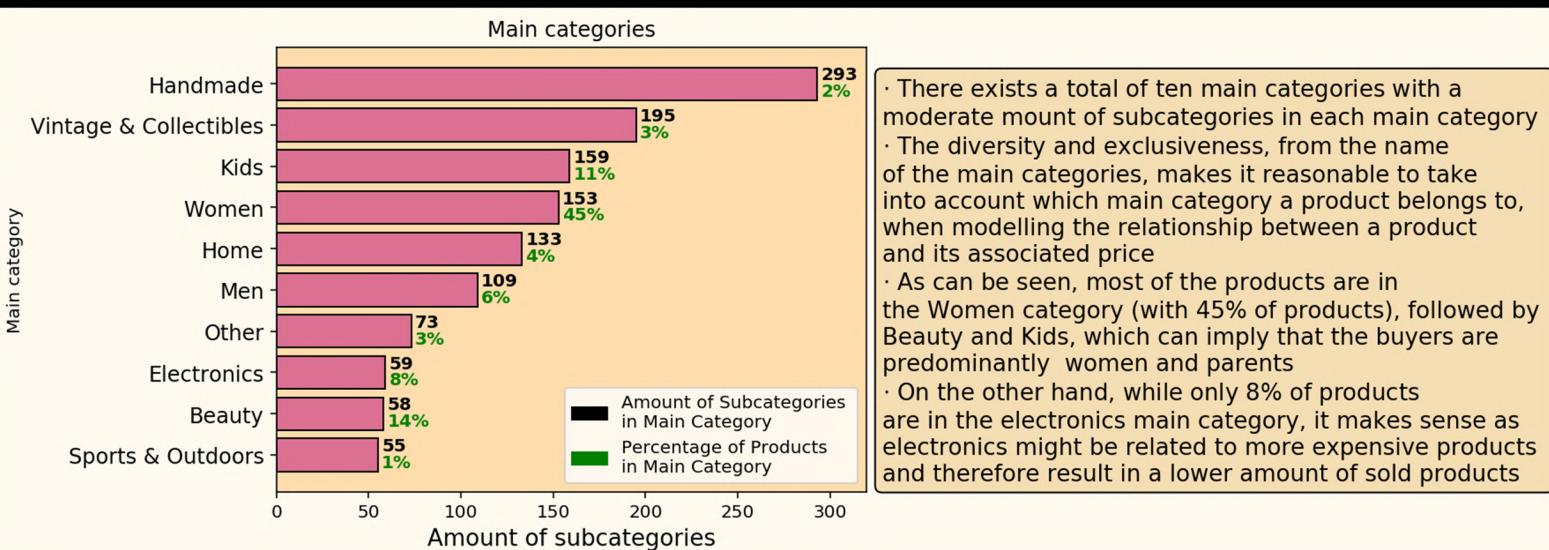
- Most products are New, followed by products in Good condition and products in Almost New condition
- A low percentage of products are either Fair or Poor, an indication that most people don't bother to post products in bad conditions
- A possible explanation is that most people tend to sell recently bought items, by e.g. regret or some other reason
- It may also indicate that buyers are mostly interested in products that are relatively new, and generally don't bother buying products with a low condition, because of e.g. less of a status symbol having low condition products

Mercari Price Analysis - Data Analysis - Categories and the main categories of the dataset

- For an idea of what type of categories each product belongs to, we consider an analysis of the category attribute -What categories exists? What are their names?

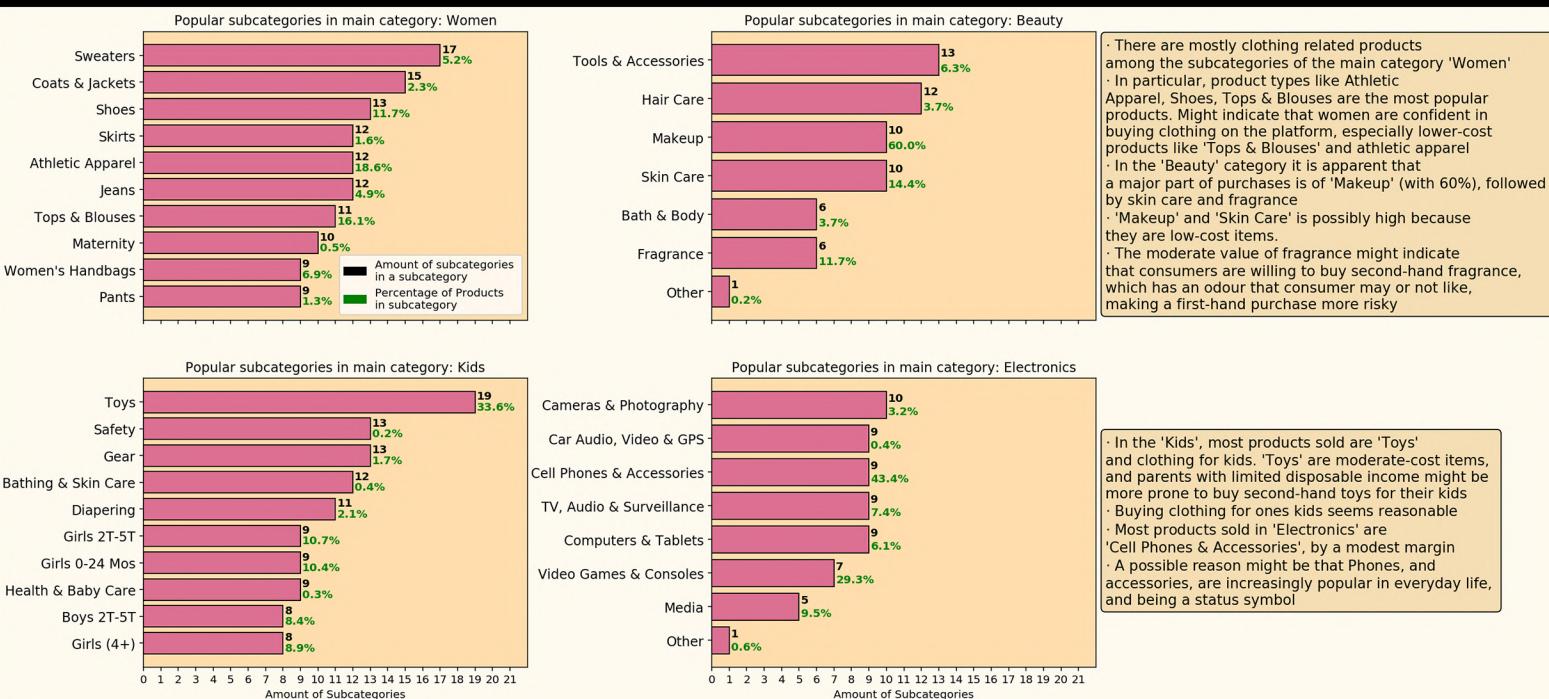


- To continue, we consider the different main categories that exists, to get a rough idea of what type of products exists



Mercari Price Analysis - Data Analysis - Subcategories in certain main categories, and a description of categories

- It is of interest to analyze what type of products exists in each main category, for this we consider some popular subcategories inside a few main categories



- The ubiquity in the amount of categories begs the question of how many categories exists at each depth

Want to quantitatively analyze the depth of categories

- Is all subcategories for a product necessary?

The amount of categories with a certain depth:

Depth	Amount of Products	Amount of Categories
3	1471819	1280
4	1330	5
5	3059	2
Total	1476208	1287

The categories with a depth of 4:

- Handmade/Housewares/Entertaining/Serving
- Men/Coats & Jackets/Flight/Bomber
- Men/Coats & Jackets/Varsity/Baseball
- Sports & Outdoors/Exercise/Dance/Ballet
- Sports & Outdoors/Outdoors/Indoor/Outdoor Games

The categories with a depth of 5:

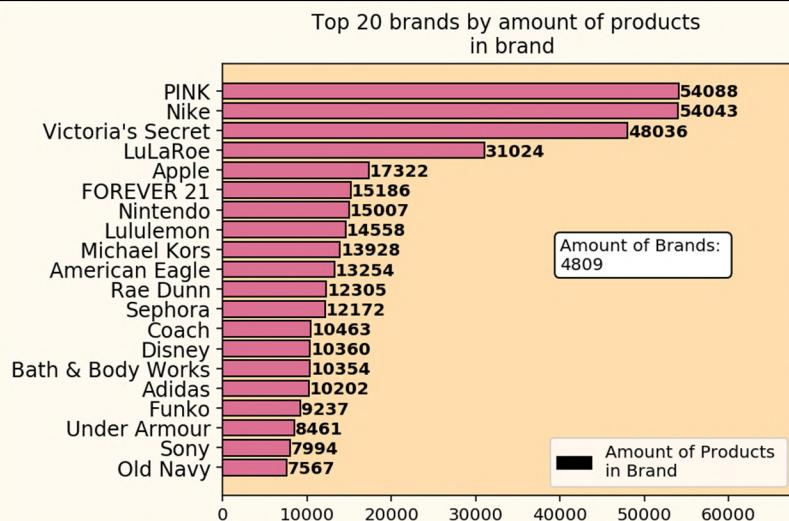
- Electronics/Computers & Tablets/iPad/Tablet/eBook Access
- Electronics/Computers & Tablets/iPad/Tablet/eBook Readers

From the structure and low quantity of the categories with a depth of 4 and 5, we can reconsider the categories as:

- Handmade/Housewares/Entertaining Serving
- Men/Coats & Jackets/Flight Bomber
- Men/Coats & Jackets /Varsity Baseball
- Sports & Outdoors/Exercise/Dance Ballet
- Sports & Outdoors/Outdoors/Indoor Outdoor Games
- Electronics/Computers & Tablets/iPad Tablet eBook Access
- Electronics/Computers & Tablets/iPad Tablet eBook Readers

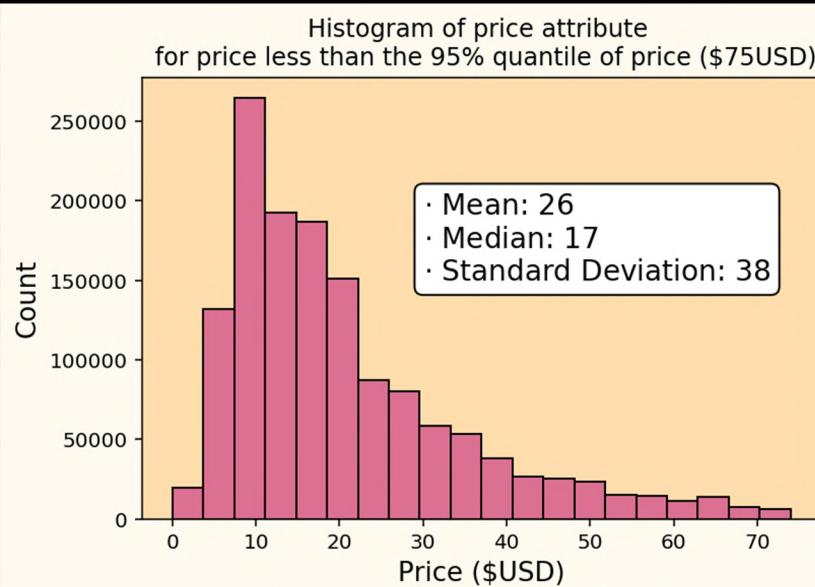
Mercari Price Analysis - Data Analysis - Brands of products, and the price of products

- Another important property of each product is what brand it belongs to - if any - as the brand would most likely have an effect on the price



- A lot of popular brands are present, like Nike, Apple, Michael Kors, Nintendo, Disney
- A total of 4809 brands indicates that it might be a good idea to analyze the average price of products in brands - to quantitatively compare the price-class of different brands, and perhaps relate certain brands to each other
- The ubiquity of popular brands is not a surprise, as there should naturally be a high demand for such products
- Another observation is how Nike is much more popular than Adidas, which is a competitor. Indicating that Nike products is to go-to sports brand for most consumers

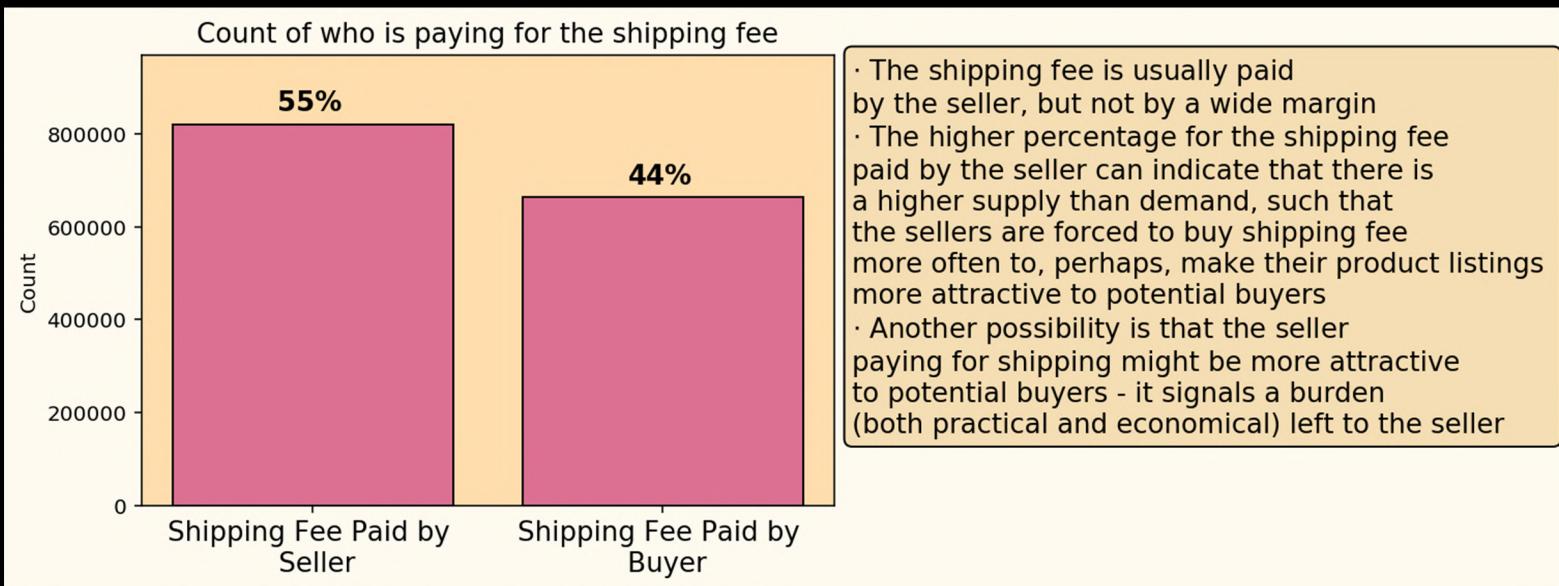
- The attribute of most interest is the price attribute, which is the price of the underlying product.



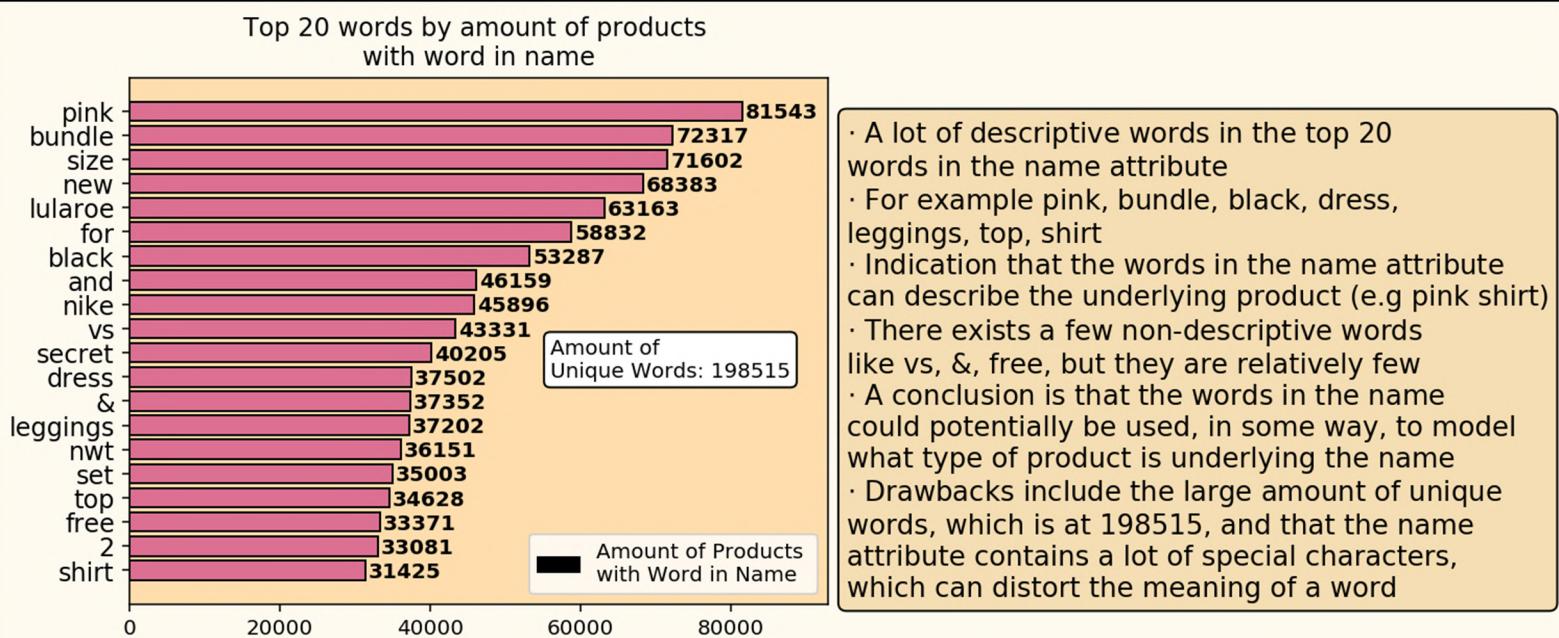
- A major part of products are bought at around 5-25 \$USD - indicates that people predominantly buy low-cost products
 - For increasing prices, the amount of products bought decreases almost exponentially - while high-cost items are bought, they are bought significantly less frequent compared to low-cost products
- The 5-25\$USD region could correspond to clothes, accessories, and other low-cost products
- It is possible that the high frequency of low-cost items bought might mean that consumers are less confident in buying high-cost products - since high-cost products might be more risky to buy second-hand rather than first-hand

Mercari Price Analysis - Data Analysis - Shipping fee and the name of each product

- For each product bought and sold on the platform, there is an associated shipping fee that is either paid by the buyer or the seller - what is its distribution?



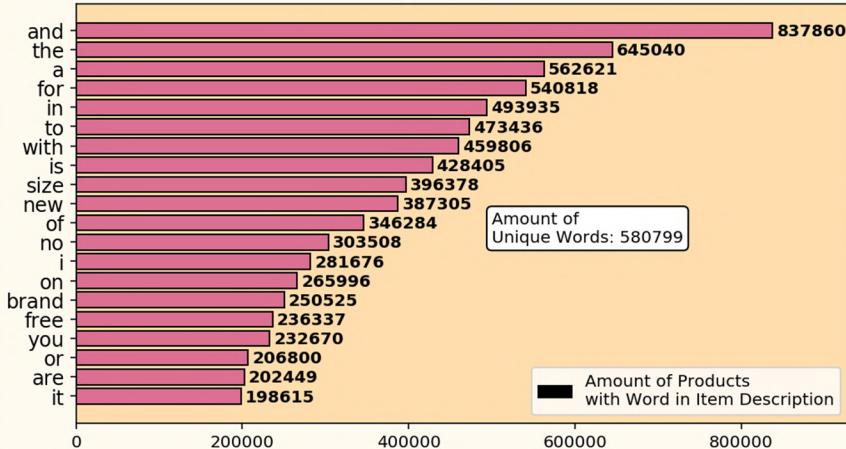
- Each product has an associated name, which is used to present the product listing to the potential buyer. Potentially, the words used in the name can be of interest in modelling the product's underlying price



Mercari Price - Data Analysis - Words in item description and the average price associated with words in the name

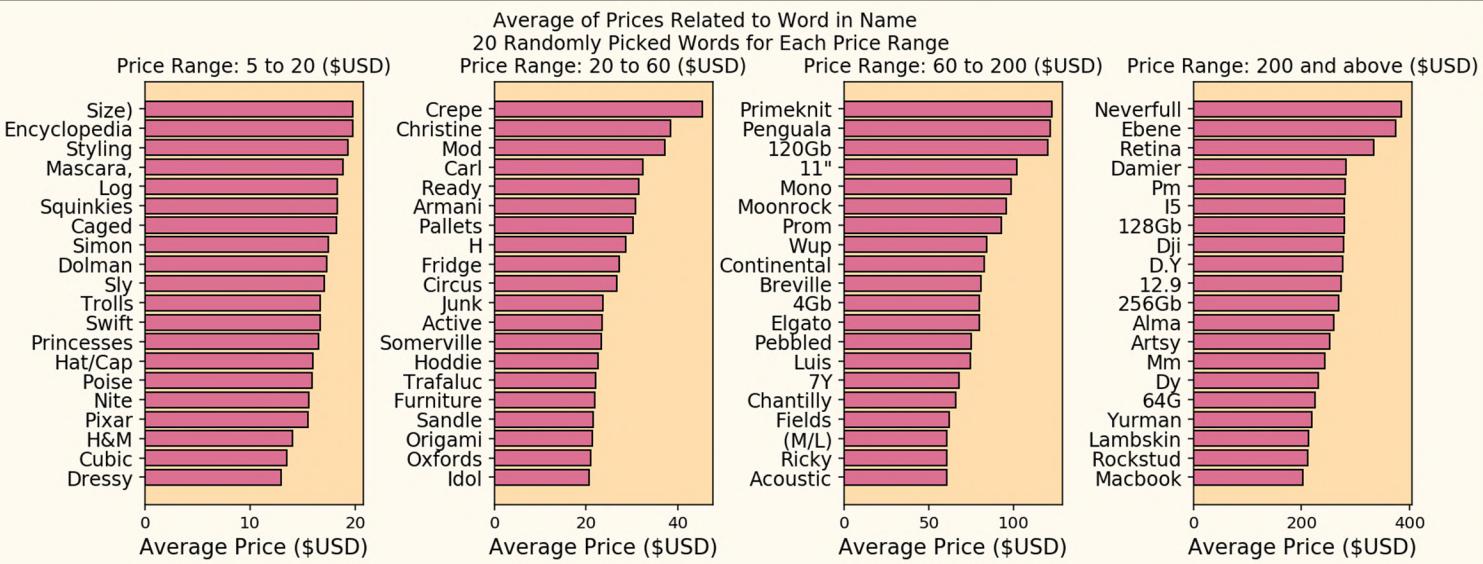
- Similarly as in the name case, the item description might contain a lot of keywords that can describe the underlying product in the listing well

Top 20 words by amount of products with word in item description



- A lot of non-descriptive words in the top 20 words in the item description
- Might imply that utilizing words in the item description might not be useful, as it contains a lot of sentence-building words, like and, the, size, brand, free, on
- In addition, a lot of words contain special characters which might distort the meaning of words
- However, all the words in the item description might collectively convey useful information, for example if the words 'size' appears with '8', it might indicate that the underlying product is a clothing piece with size 8. This type of inference should be able to convey some information on the possible price of the underlying item
- Hence, the words collectively might convey a lot of useful information

- In addition to what words exists in the name of products, there is an interest in correlating those words with prices of underlying products

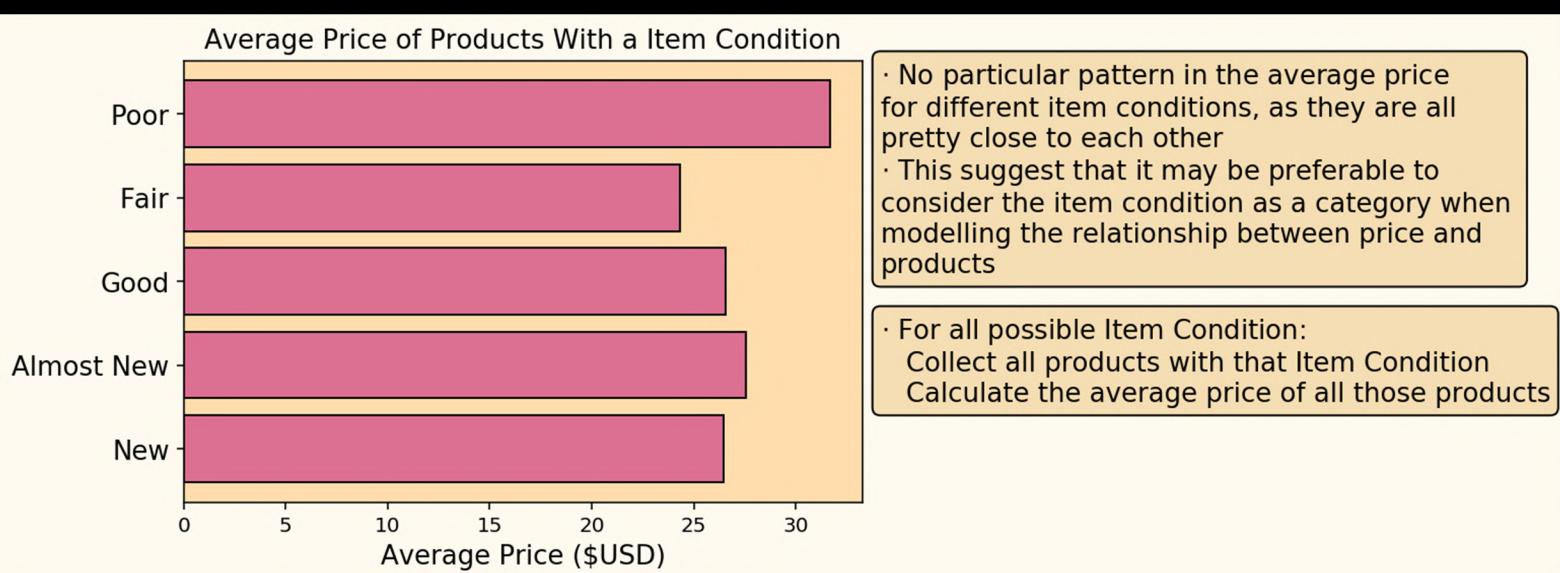


- A lot of descriptive words in all price ranges, including words like Trolls, Princesses, Armani, Hoddie, Oxfords, Acoustic, 120Gb, 256Gb, Macbook, Lambskin, Damier
- Some words give an almost full description of the underlying product, e.g. Macbook, while some words give an important aspect, e.g. Lambskin
- Hence, it is concluded that the words, especially together, in the name attribute of a product can give potentially useful information of a product's underlying price

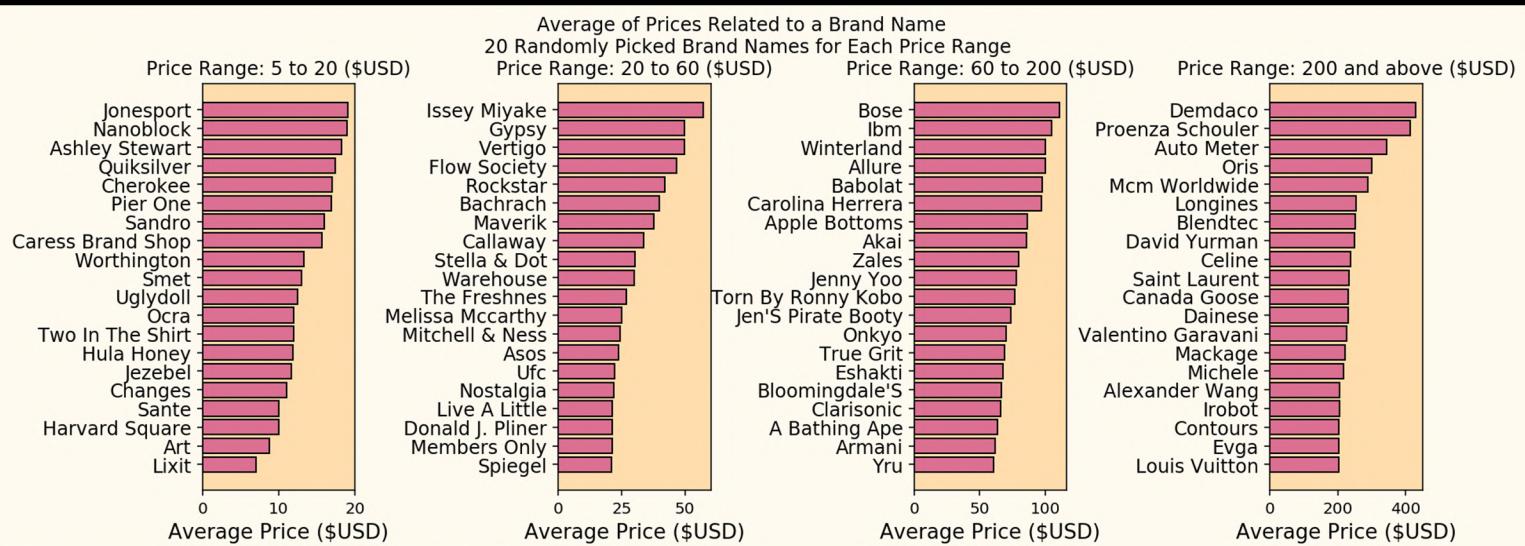
- Construction of the Average Price: For each possible Word in the name attribute: Collect all products which has that Word in its name attribute Calculate the average price from all the prices of those products

Mercari Price Analysis - Data Analysis - Prices associated with item condition of products, and prices of products in different brands

- A minor interest is to consider the prices of products with different item conditions - are less quality items generally cheaper?



- A major factor of the price will obviously be which brand the underlying product belongs to - especially what type of price products in that brand usually sell for, in a general sense

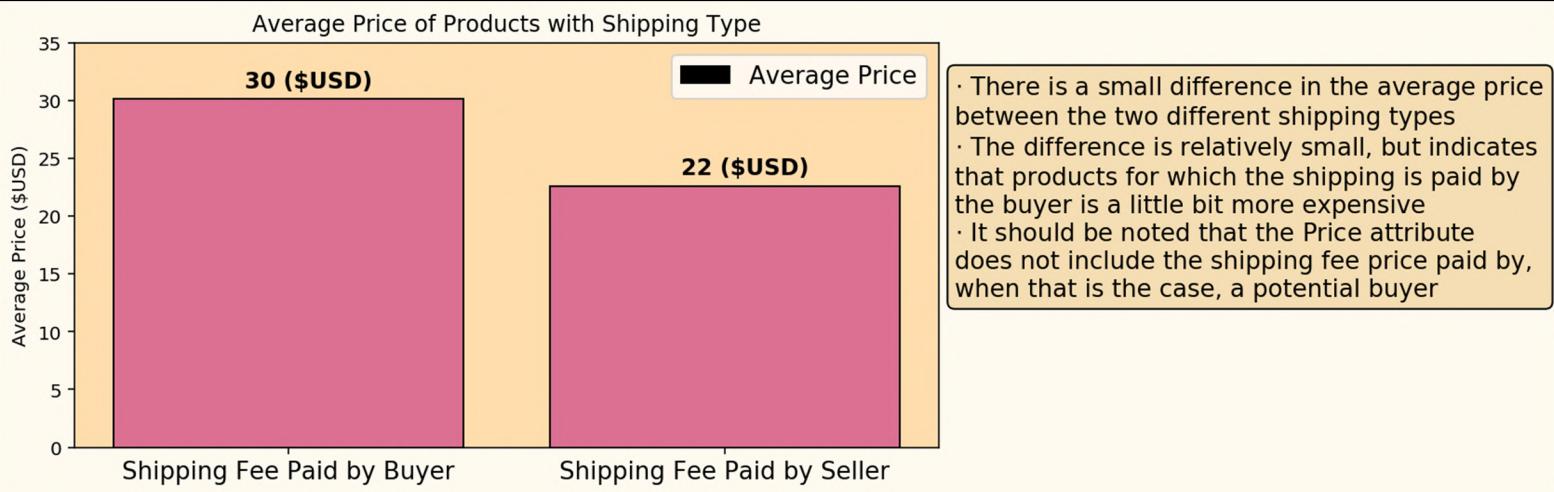


- A lot of known brands in all price ranges, including brands like Asos, David Yurman, Louis Vuitton, Saint Laurent, Alexander Wang
 - The average price of products associated with a brand can give a rough idea of the price-class associated with the brand - some may be expensive like Louis Vuitton while others may be more affordable like Pier One
 - In addition, certain brands will have a similar price-class, and will probably sell products for similar prices

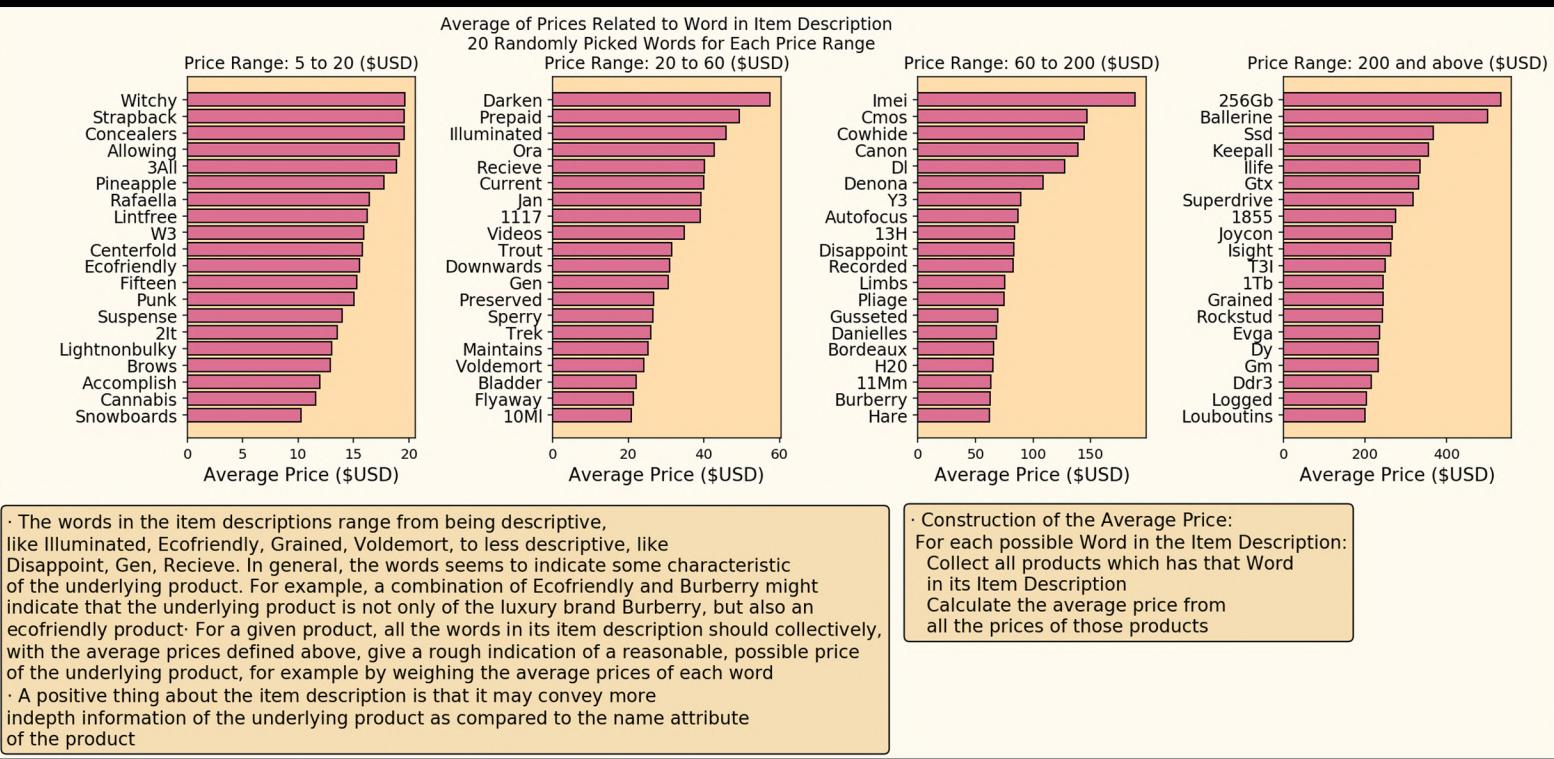
- Construction of the Average Price:
For each possible Brand Name:
Collect all products associated with that Brand Name
Calculate the average price of all those products

Mercari Price Analysis - Data Analysis - Shipping fee with price, and the price associated with words in the item description

- Further considering the shipping fee of items, there might be a relation between who pays the shipping fee and the average price of products

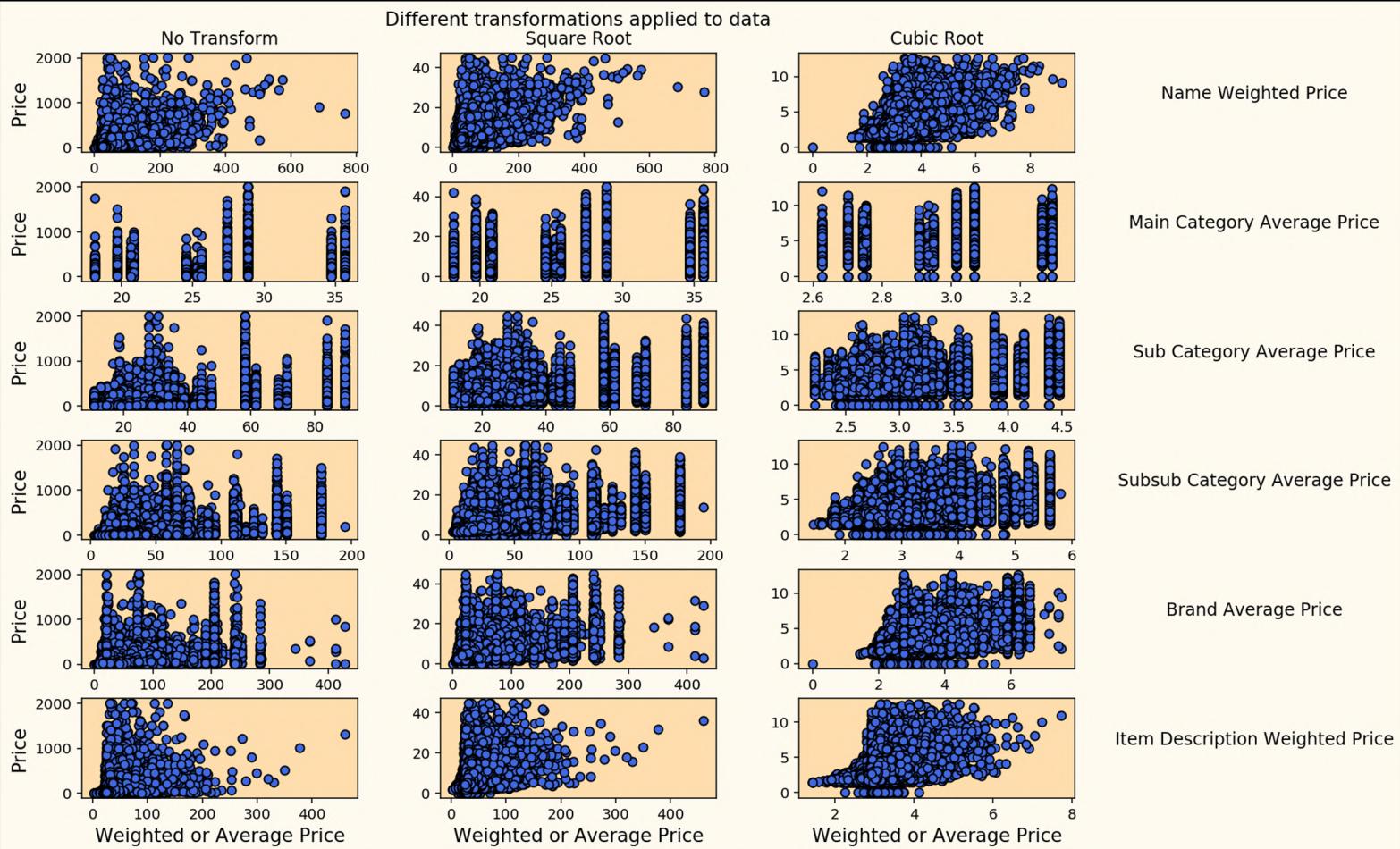


- The item description might convey important information of the underlying product - hence, the price associated with words in the item description is of interest



Mercari Price Analysis - Prediction Analysis - Transformation of variables, correlation among variables

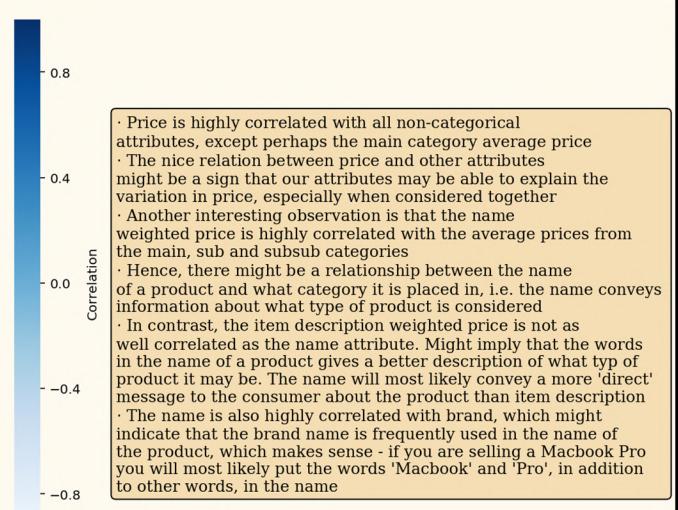
- Before building a linear model, there might be some transformations of variables that may help with relationships



- The cubic root transformation makes the different attributes a bit more linearly related with the price attribute
- On the other hand, one can see that the large amount of data points (above 1 million) makes the relationships really crowded

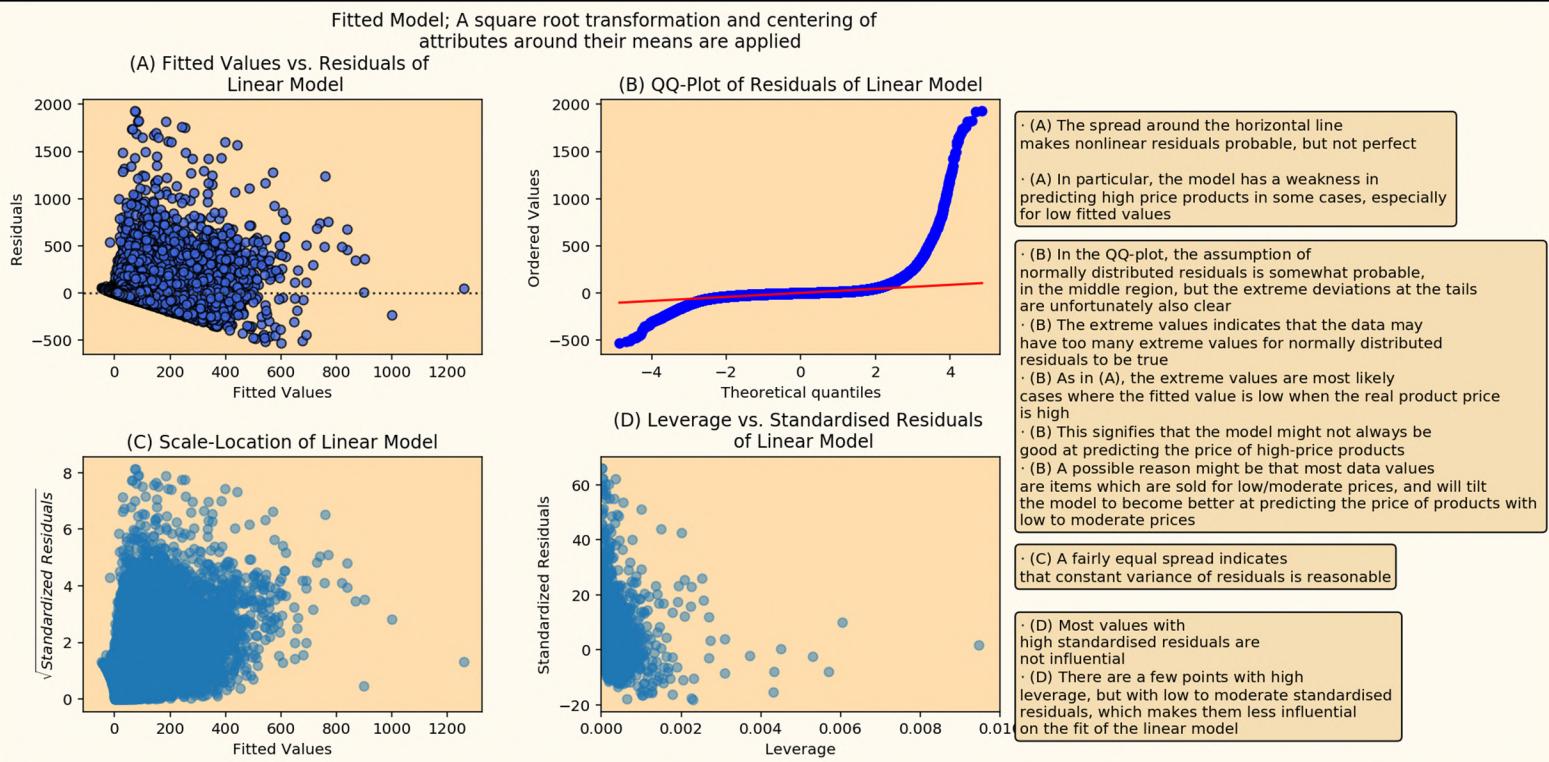
- A key idea in linear models is to analyze the correlation among all variables considered, especially the correlations with respect to the price variable

Correlation Between Different Attributes														
Name Weighted Price	1	-0.07	0.0035	0.051	0.061	0.037	0.24	0.47	0.55	0.6	0.53	0.11	-0.11	0.64
Item Condition Id 1.0	-0.07	1	-0.51	-0.56	-0.13	-0.035	-0.19	-0.11	-0.14	-0.07	-0.082	-0.2	0.2	-0.0035
Item Condition Id 2.0	0.0035	-0.51	1	-0.37	-0.087	-0.023	0.096	0.028	0.026	0.0061	0.015	0.077	-0.077	0.015
Item Condition Id 3.0	0.051	-0.56	-0.37	1	-0.095	-0.026	0.097	0.071	0.093	0.045	0.057	0.13	-0.13	-0.0061
Item Condition Id 4.0	0.061	-0.13	-0.087	-0.095	1	-0.006	0.035	0.068	0.086	0.067	0.048	0.037	-0.037	-0.016
Item Condition Id 5.0	0.037	-0.035	-0.023	-0.026	-0.006	1	0.033	0.033	0.07	0.046	0.037	0.007	-0.007	0.0068
Main Category Average Price	0.24	-0.19	0.096	0.097	0.035	0.033	1	0.5	0.34	0.27	0.22	0.071	-0.071	0.17
Sub Category Average Price	0.47	-0.11	0.028	0.071	0.068	0.033	0.5	1	0.72	0.39	0.39	0.11	-0.11	0.34
Subsub Category Average Price	0.55	-0.14	0.026	0.093	0.086	0.07	0.34	0.72	1	0.4	0.45	0.13	-0.13	0.48
Brand Average Price	0.6	-0.07	0.0061	0.045	0.067	0.046	0.27	0.39	0.4	1	0.41	0.068	-0.068	0.48
Item Description Weighted Price	0.53	-0.082	0.015	0.057	0.048	0.037	0.22	0.39	0.45	0.41	1	0.09	-0.09	0.49
Shipping 0.0	0.11	-0.2	0.077	0.13	0.037	0.007	0.071	0.11	0.13	0.068	0.09	1	-1	0.2
Shipping 1.0	-0.11	0.2	-0.077	-0.13	-0.037	-0.007	-0.071	-0.11	-0.13	-0.068	-0.09	-1	1	-0.2
Price	0.64	-0.0035	0.015	-0.0061	-0.016	0.0068	0.17	0.34	0.48	0.48	0.49	0.2	-0.2	1



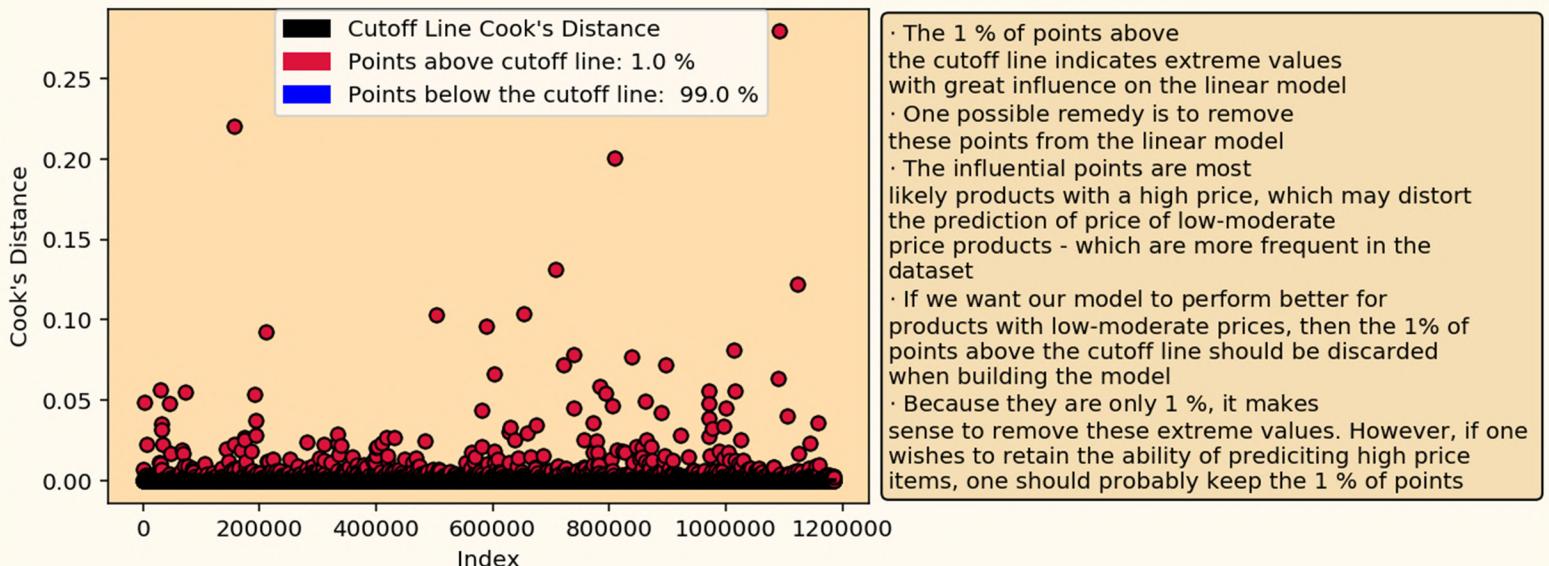
Mercari Price Analysis - Prediction Analysis - First linear model; Diagnostics and Cook's Distance

- A first linear model is applied to the problem, with the price attribute used as a response variable. To analyze the appropriateness of the model, the idea is to consider diagnostic plots



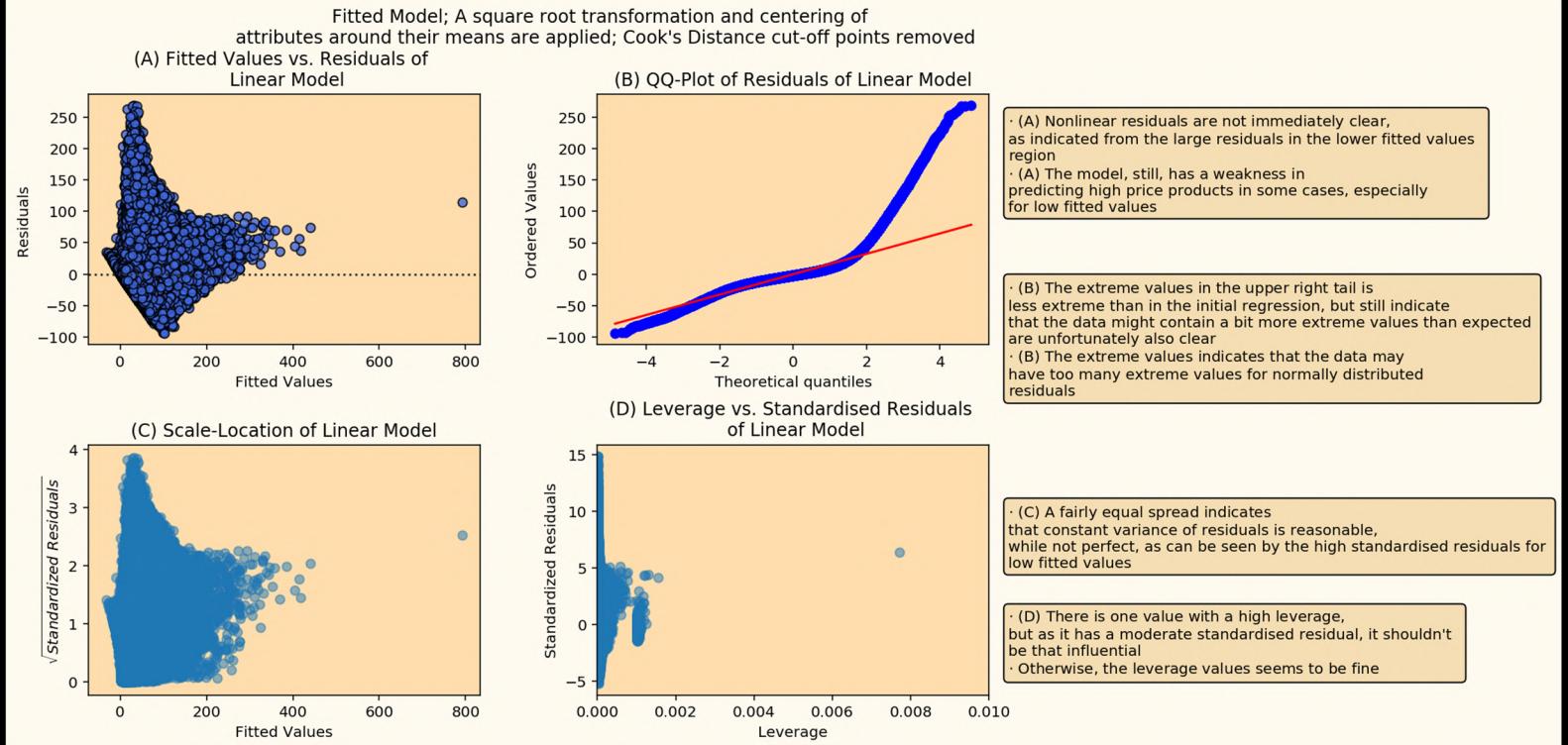
- Further describing the linear model, especially illustrating data points with influence on the fit of the model, the Cook's Distance can be considered

Cook's Distance of Fitted Values

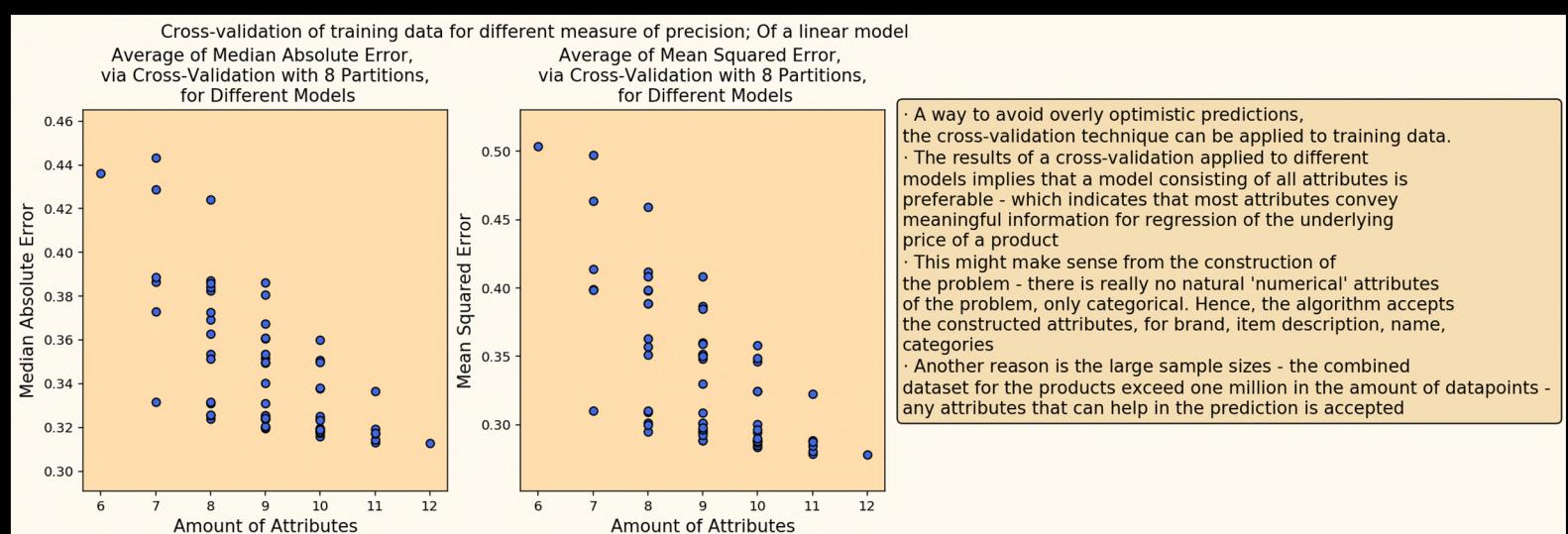


Mercari Price Analysis - Prediction Analysis - Linear model after Cook's Distance analysis, and different performance measures to choose an optimal model

- After removing datapoints based on Cook's distance, another linear model is built and its properties are analyzed



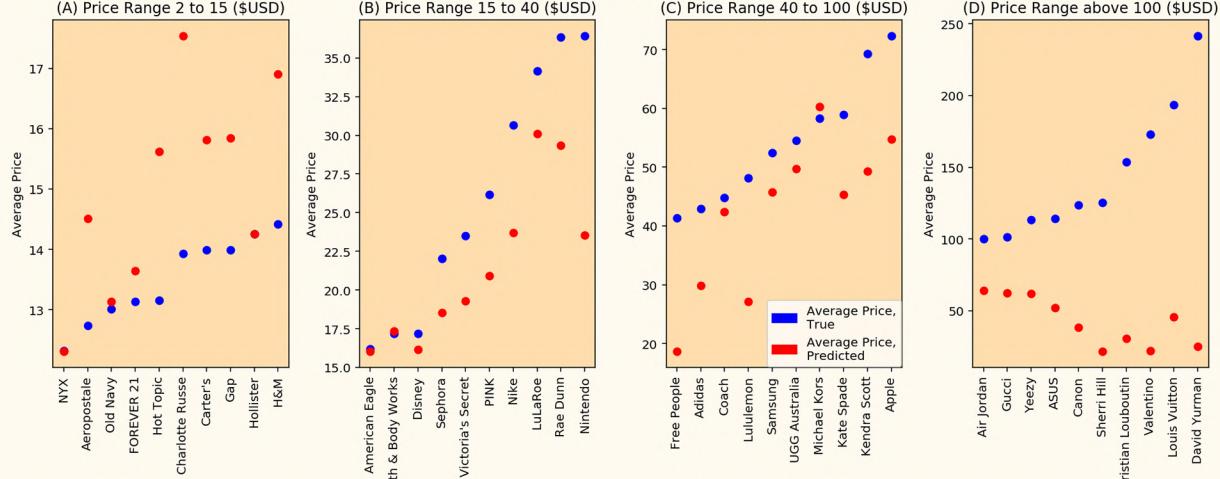
- The aim is to find an optimal model. For this, we consider all combinations of attributes (where all binary attributes are kept in the model), and for each combination we evaluate two performance measures. Lastly, based on the performances, we pick a final, optimal model



Mercari Price Analysis - Prediction Analysis - Optimal model in predicting prices of products in brands, and predicting products in different main categories

- With our final model, the aim is to evaluate the model on datasets not used in designing the model. The idea is measure how well the model performs in practice. For a starter, we consider predictions on products in brands - to see how well the model performs when different brands of products is considered

Predicted and True Average Price of Products in Brands; Ten Brands in Various Price Ranges;
Predictions on a Validation Set (Independent of Training Set)



(A) The predictions in the low range performs remarkably well, where all predictions are just a few dollars away from the true values.
 (A) Indicates that our model is in general good for predicting product prices for low-price range brands like H&M, Old Navy and Gap
 (A) A possible reason might be that a lot of products exists in the low-priced range, which makes our model really efficient in predicting such products

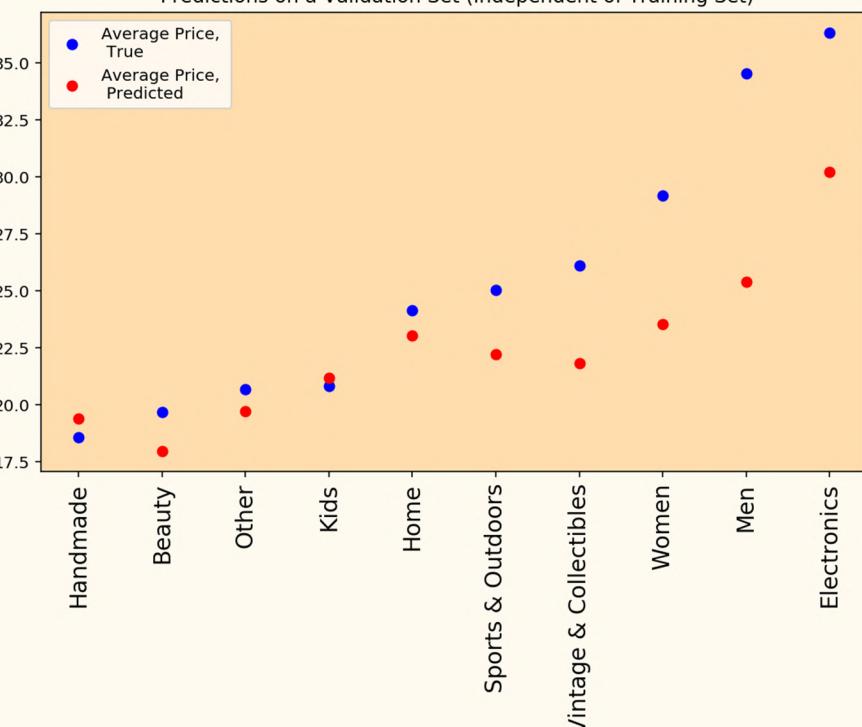
(B) Similarly as in the low-price range, the prediction of price in moderate price range is really well too - again with just a few dollars of margin
 (B) Hence, our model should perform well predicting prices of products like Nike, Disney and American Eagle
 (B) The well performance might come, as in (A), from the ubiquity of moderate priced products

(C) The performance in predicting high priced products differs among brands, as can be seen
 (C) While a few brands can be predicted well, like Coach, Samsung, Michael Kors, the model has a difficulty predicting e.g. Free People, Lululemon
 (C) However, in general, the predictions are not too bad. One could argue that a margin of 10-15 (\$USD) among high-priced products isn't that much, especially for brands like Apple

(D) Contrary to the other price ranges, the predictions for very-high priced products isn't generally well, as can be seen from the increased margin as the true price increases
 (D) Most likely, it is due to two factors. Firstly, in our model a portion of training data points corresponding to high cook's distance points were removed, which in the process have decreased our model's ability to predict prices for very-high priced products. Secondly, there is most likely not as many high-priced products in our training data, as opposed to low-, moderate-priced products, which most likely have made our model biased towards predicting the correct price of low- and moderate-priced products

- Further, there is an interest in evaluating our optimal model on products from different main categories.

Predicted and True Average Prices of Products in Main Category;
Predictions on a Validation Set (Independent of Training Set)



The predicted values seem to be quite good, in general, for most of the main categories
 However, there seems to be some deviations when predicting products in the Women, Men and Electronics main categories. Possibly because products in these categories range from low-priced to very-high priced brands. A major portion of very-high priced brands will have a negative effect on the predictions, as discussed in the predictions of products in brands

Mercari Price Analysis - Prediction Analysis - Optimal model in predicting products in different categories

- Lastly, there is an interest in picking a few particular categories of products and evaluate our optimal model on these categories. The idea is that these categories corresponds to a diverse set of products, and this will show how the model performs on a diverse set of categories

