

# ML Special Assignment

- Munmun Sangtani, 20BEC073
- Siddharth Agarwal, 20BEC119

# What we have developed?

## E-Commerce Product Price Prediction

We aim to develop a pricing algorithm that automatically suggests optimal prices for products listed on online platforms.

# Why we have developed this?

By providing users with suggested prices based on product attributes and user-inputted text descriptions, enhance the buying and selling experience can be enhanced and the volume of transactions can be increased.

# How we have developed this?

we have presented a comparative study of two machine learning algorithms, LightGBM and Ridge Regression, for e-commerce price prediction tasks. Additionally, we have evaluated their performance using the RMSLE metric, which is particularly useful for predicting logarithmic values.

# Flow of Program:

Importing dataset

Preparing the Corpus for Analysis

Explore Training Set

Create three new features from Categories (Main, Sub1, Sub2)

Create CSR\_Matrix & Merge the Sparse Matrices

Data Cross validation - Holdout Method

Train with LGBM

Train with Ridge Regression

Predict on test set

# Dataset

Train-id / test-id	id of the listing
Name	the title of listing.
condition-id	the condition of the items by seller
Category-name	category of listing
Brand-name	brand of listing
Price	the price for the item sold. Our Target variable we want to predict, which is in USD.
Shipping	1 if shipping fee paid by seller and 0 by buyer
Item-description	the full description of item.

# Algorithms Used

# Algorithms Used

- **LGBM** : It is a gradient boosting framework that uses a tree-based learning algorithm and is designed to be efficient, scalable, and accurate for large-scale datasets.
- rmsle: 0.6571443747117787
- **Ridge Regression**: It is a regularization technique that adds a penalty term to the regression model's cost function to reduce the impact of multicollinearity on the model's coefficients.
- rmsle: 0.47049933995488696

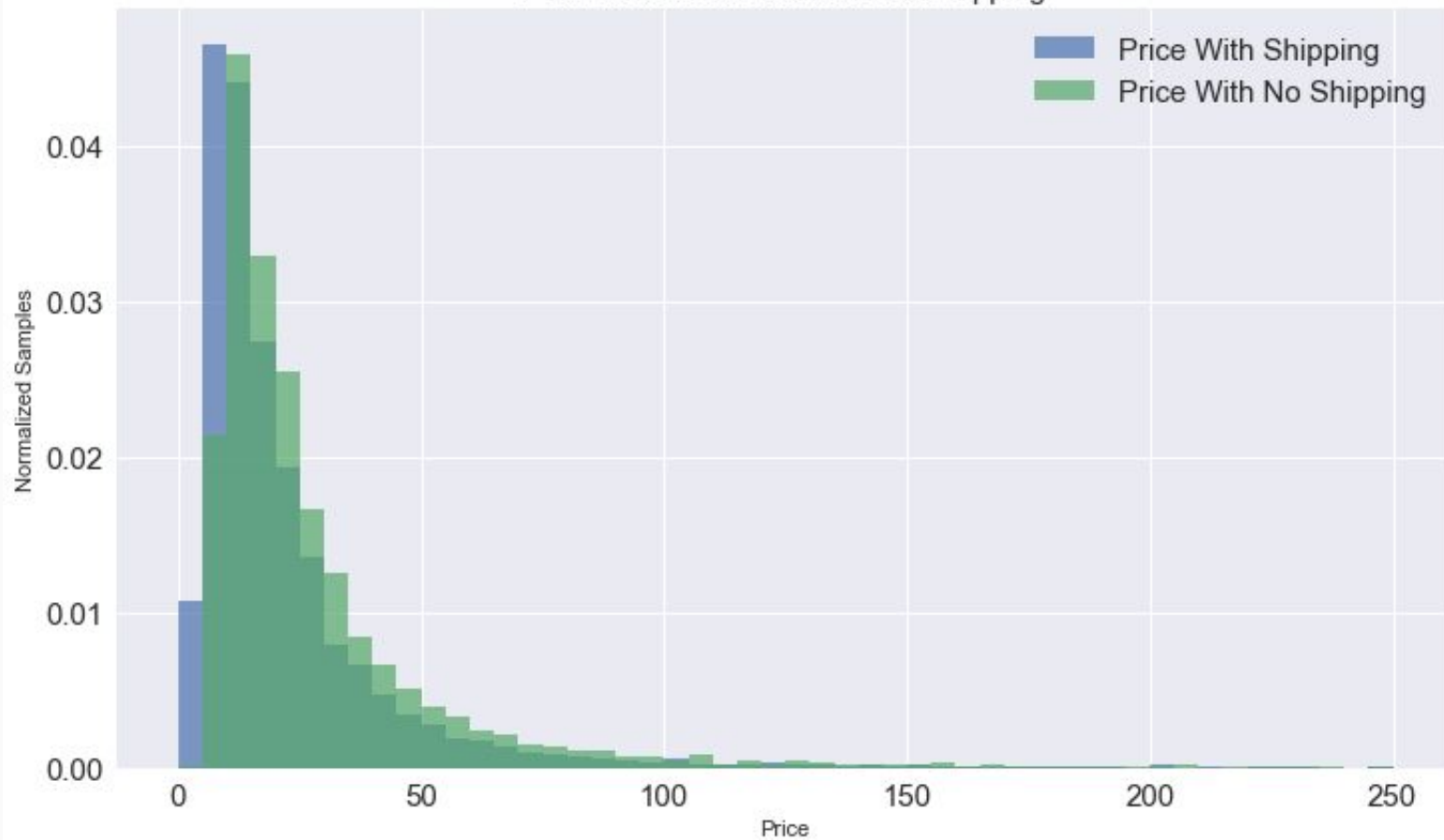


# Results and Analysis

# Price Distribution

- The mean price in the dataset is 26 Dollars
- The median price in the dataset is 17 Dollars
- The max price in the dataset is 2000 Dollars
- Due to the skewed dataset, the median price is a more reliable price to gauge off of.

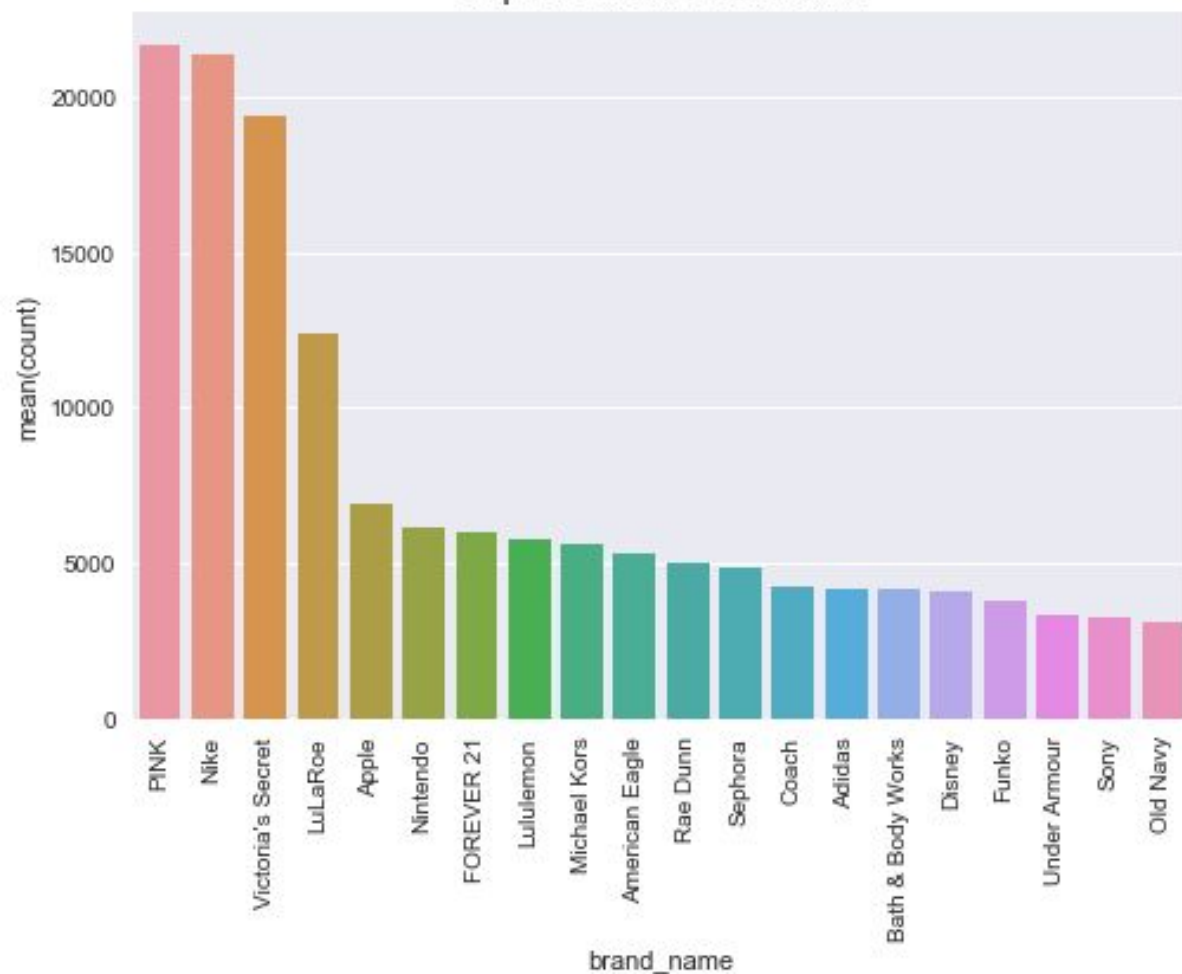
Price Distrubtion With/Without Shipping



# Top 20 brand distribution

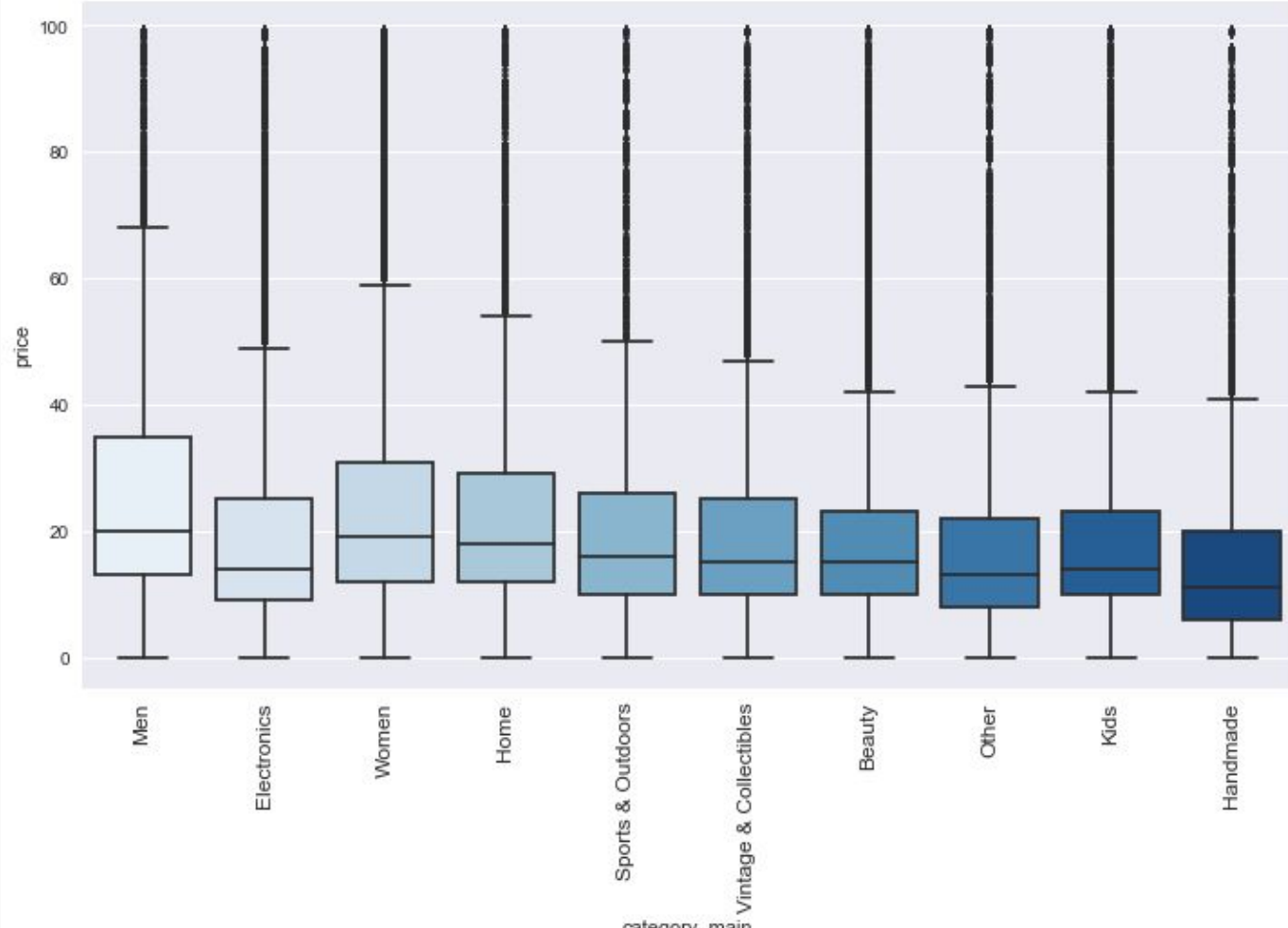
Majority of the top brands are clothing brands and electronics. PINK and Victoria Secret are among the top 3 brands and are typically towards female customers.

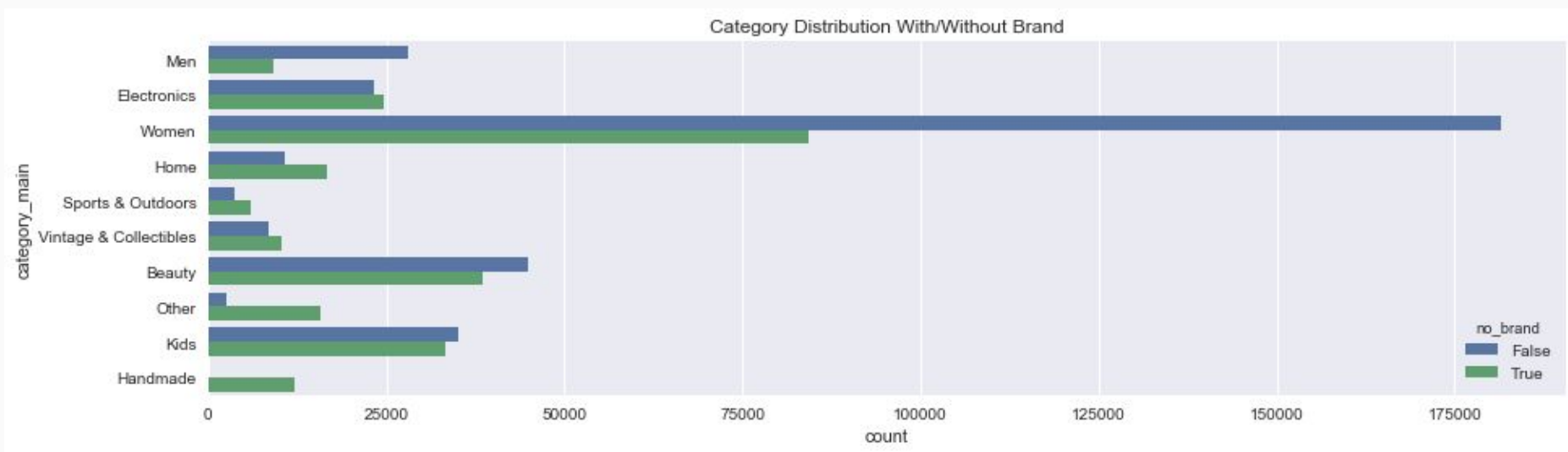
Top 20 Brand Distribution



# Main Category

It is seen that majority of the distribution is taken by women and beauty. They take 56% of the distribution. The prices are evenly distributed across all categories. The Men category the only one that averages out the most.







# Effect of word count

Effect of word count of item description on prices: It was observed that from about 0-300 words, there is a positive linear relationship. After that a gradual negative relationship was seen. It drops at about the 1000 word point. Overall, the word count did not have any significant role in deciding the price of an item

