# PREDICTING THE PRICE OF SKINCARE PRODUCTS USING REGRESSION MODEL

JING LEE

APRIL 20, 2022

METIS PROJECT 2

WEB SCRAPING AND REGRESSION

# Introduction

- **Motivation**
  - Skincare is a multibillion-dollar growing industry
  - The price for 1 oz cream range from $10 to $200+
  - What determine the pricing of products?
  - Budget brands vs luxury brands?
- **Objective**
  - Investigation the features that contribute to the pricing of products
  - Use a regression model to predict price of skincare products
- **Goals**
  - Consumers can decide what products to buy within their budget
  - New skincare companies can understand how products are priced

# What goes into the pricing of products?



- $123 per 1 oz
- Parent company: L'oreal
- Luxury brand
- Rating: 4.6 stars
- Reviews: 3,919



- $17 for 1 oz
- Parent company: L'oreal
- Budget brand
- Rating 4.12 stars
- Reviews: 112

# Methodology: data collection

- Web scraping data (~1847 products)
  - website: www.ulta.com
  - Major retailer of beauty products ranging from luxury to affordable brands

# Methodology: Tools

- Tools:
  - Requests/ BeautifulSoup, pandas, matplotlib, seaborn, statsmodels, sklearn
- Models:
  - Linear regression without regularization
  - Linear regression with regularization
    - Lasso
    - Ridge
    - Elastic net
  - Extensive feature engineering to improve model
- Metrics:
  - $R^2$, mean absolute error (MAE)

# Model selection – cross validation

- Linear regression without regularization
  - Predictors: $10 \rightarrow 17$ and 1501 datapoints
- Mean train $R^2$ : 0.409 +/- 0.063
- Test $R^2$ : 0.413
- MAE: $13.66

# Predictors in the model

- Predictors correlation with price (predictor)
  - len_ingredients  0.44
  - price_per_vol 0.44



Skincare products pricing predictors

# Other findings

- Most brand come have the same  7 major parent companies
- High ratings and low number of reviews and vice verse



Distribution of number of reviews

Distribution of ratings

# Conclusions

- Predicting price with linear regression and the current features is not a good model

- Model has low prediction accuracy of ~ 41 %

- Price prediction is +/- $13.66 is high considering average products cost ~$35

- Missing important features such as investment in marketing which could heavily influence pricing and purchase power

- There is no clear predictor for products pricing

# Future work

- Acquire data from other products sources
  - Sephora: luxury brands
  - Target: affordable brands
- In depth study of the actual ingredients in the product to determine pricing
  - use of classification model to improve prediction
- A feature to account for marketing influence

# Thank you!!

Questions?

# Appendix

```
simple regression scores:
train r2 [0.39054358 0.41968104 0.37822006 0.52167151 0.33326757]
mean cv r2 0.409 +/- 0.063
test r2: 0.413
MAE +- $: 13.66

ridge regression scores:
train r2 [0.39237975 0.43186739 0.36834842 0.51323988 0.3546437 ]
mean cv r2 0.412 +/- 0.057
test r2: 0.409
MAE +- $: 13.78

lasso regression scores:
train r2 [0.39058034 0.42079321 0.37841136 0.52173008 0.33252066]
mean cv r2 0.409 +/- 0.063
test r2: 0.413
MAE +- $: 13.66

elastic regression scores:
train r2 [0.39187371 0.43534027 0.36727406 0.5114768  0.35664375]
mean cv r2 0.413 +/- 0.056
test r2: 0.408
MAE +- $: 13.81
```
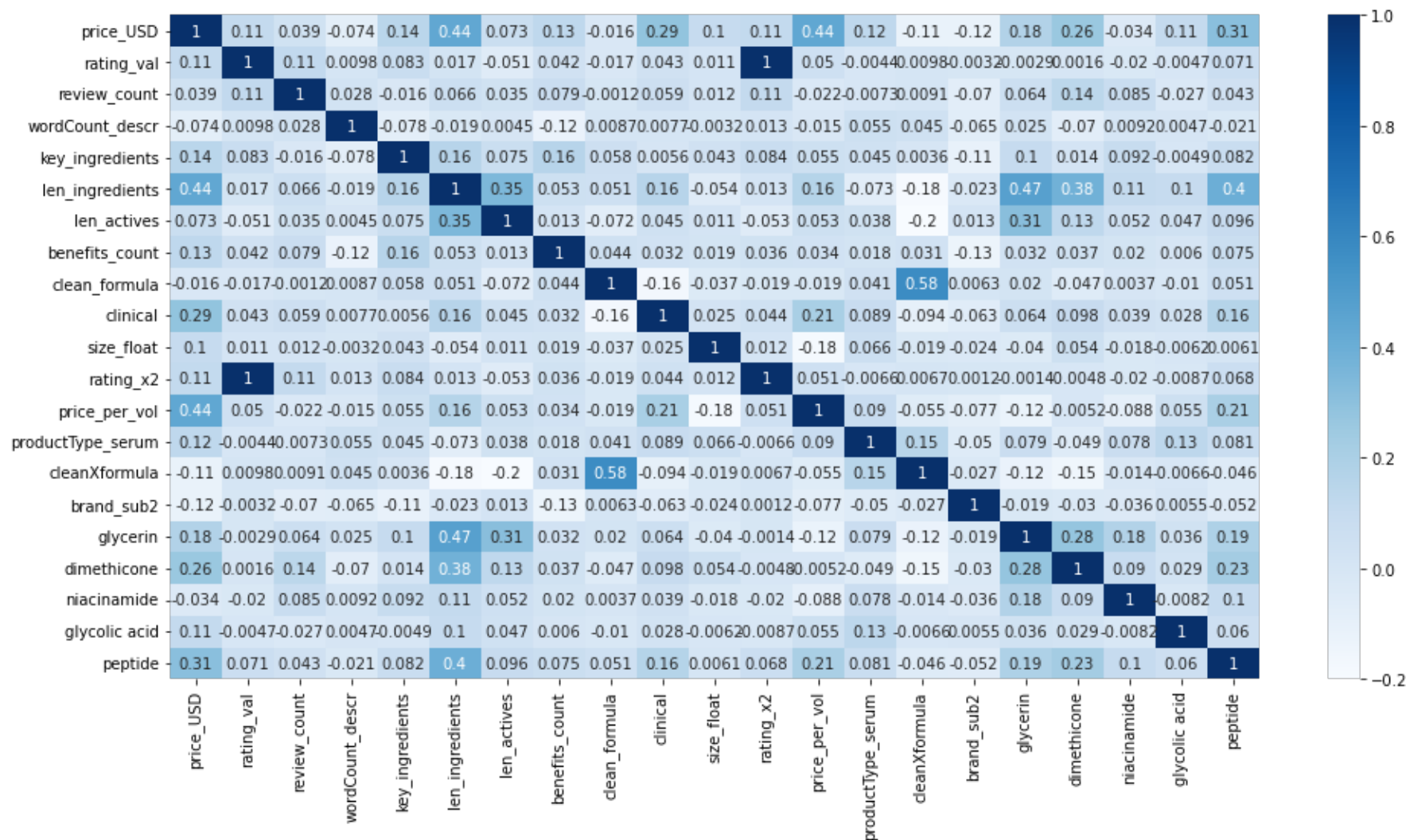
# Heatmap

|  | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| const | 15.0415 | 5.738 | 2.621 | 0.009 | 3.786 | 26.297 |
| wordCount_descr | -0.1526 | 0.057 | -2.680 | 0.007 | -0.264 | -0.041 |
| key_ingredients | 0.4733 | 0.236 | 2.002 | 0.045 | 0.010 | 0.937 |
| len_ingredients | 0.4791 | 0.041 | 11.610 | 0.000 | 0.398 | 0.560 |
| len_actives | -3.3086 | 0.671 | -4.932 | 0.000 | -4.625 | -1.993 |
| benefits_count | 0.4320 | 0.138 | 3.123 | 0.002 | 0.161 | 0.703 |
| clinical | 1.4984 | 0.250 | 5.990 | 0.000 | 1.008 | 1.989 |
| size_float | 0.7687 | 0.093 | 8.255 | 0.000 | 0.586 | 0.951 |
| rating_x2 | 0.5725 | 0.164 | 3.499 | 0.000 | 0.252 | 0.893 |
| price_per_vol | 0.2715 | 0.016 | 16.936 | 0.000 | 0.240 | 0.303 |
| productType_serum | 5.2029 | 1.112 | 4.678 | 0.000 | 3.021 | 7.384 |
| cleanXformula | -0.3645 | 0.205 | -1.774 | 0.076 | -0.767 | 0.039 |
| brand_sub2 | -6.9184 | 2.737 | -2.528 | 0.012 | -12.287 | -1.550 |
| glycerin | 3.6681 | 1.345 | 2.728 | 0.006 | 1.031 | 6.306 |
| dimethicone | 5.4102 | 1.100 | 4.917 | 0.000 | 3.252 | 7.569 |
| niacinamide | -5.1606 | 1.502 | -3.436 | 0.001 | -8.107 | -2.214 |
| glycolic acid | 4.0026 | 2.174 | 1.841 | 0.066 | -0.262 | 8.267 |
| peptide | 2.8370 | 1.383 | 2.051 | 0.040 | 0.124 | 5.550 |

| | | | |
|---|---|---|---|
| Omnibus: | 162.332 | Durbin-Watson: | 1.982 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 1318.181 |
| Skew: | -0.072 | Prob(JB): | 5.76e-287 |
| Kurtosis: | 7.589 | Cond. No. | 730. |

13