# Sephora

What Drives Product Recommendations?

IRONHACK FINAL PROJECT

Elham Allahbakhshi

# Project Overview



- **Product name**
- **Brand name**
- **Price**
- **Discount price**
- **Sephora edition**
- **Rating**
- **Review text**
- **Review title**
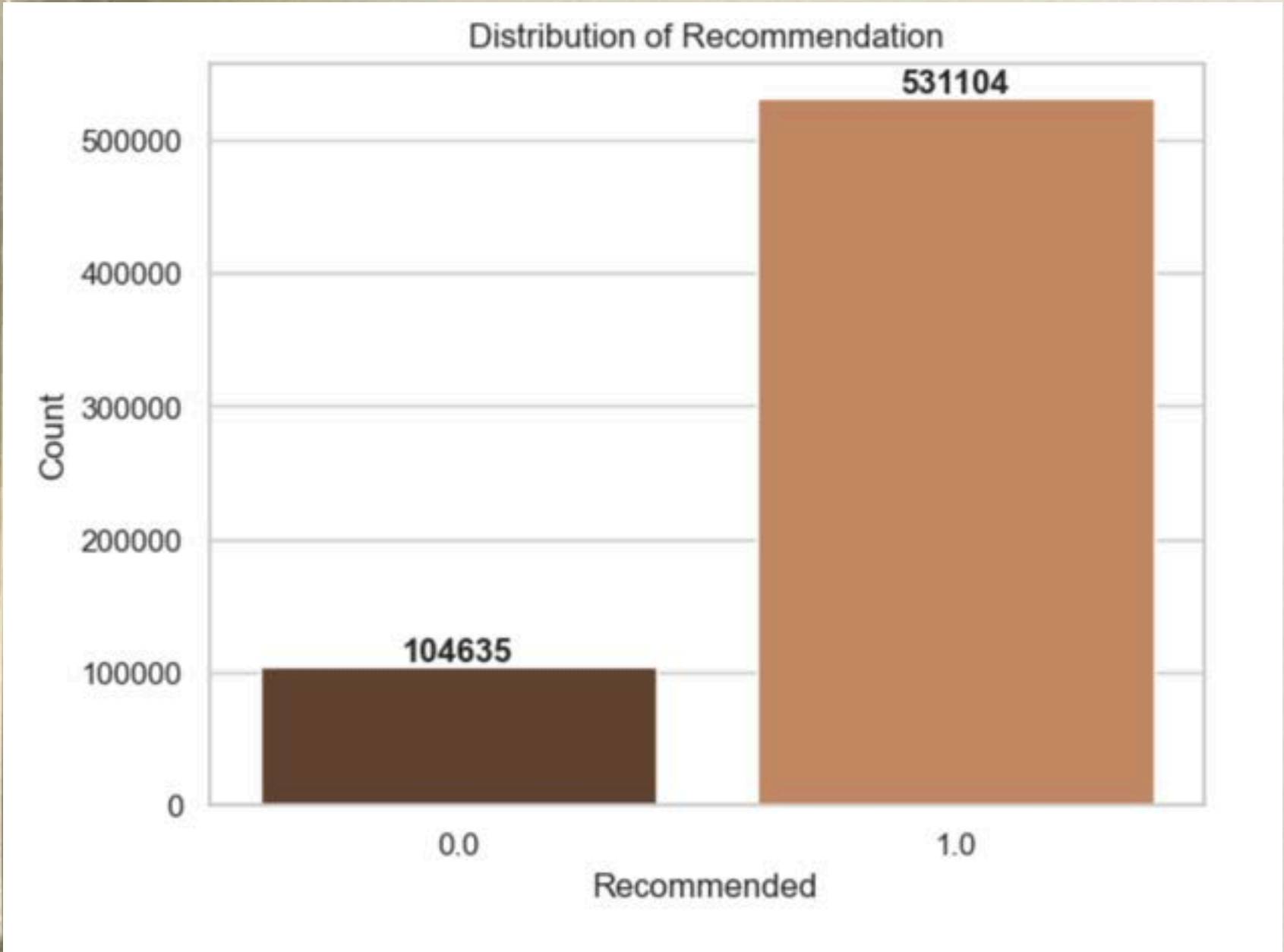- **Skin_tone**
- **Eye_color**
- **Skin_type**
- **Category**
- **Is_recommended**

Project objective: Identifying factors that influenced consumers recommending a product

# Distribution of Recommendation



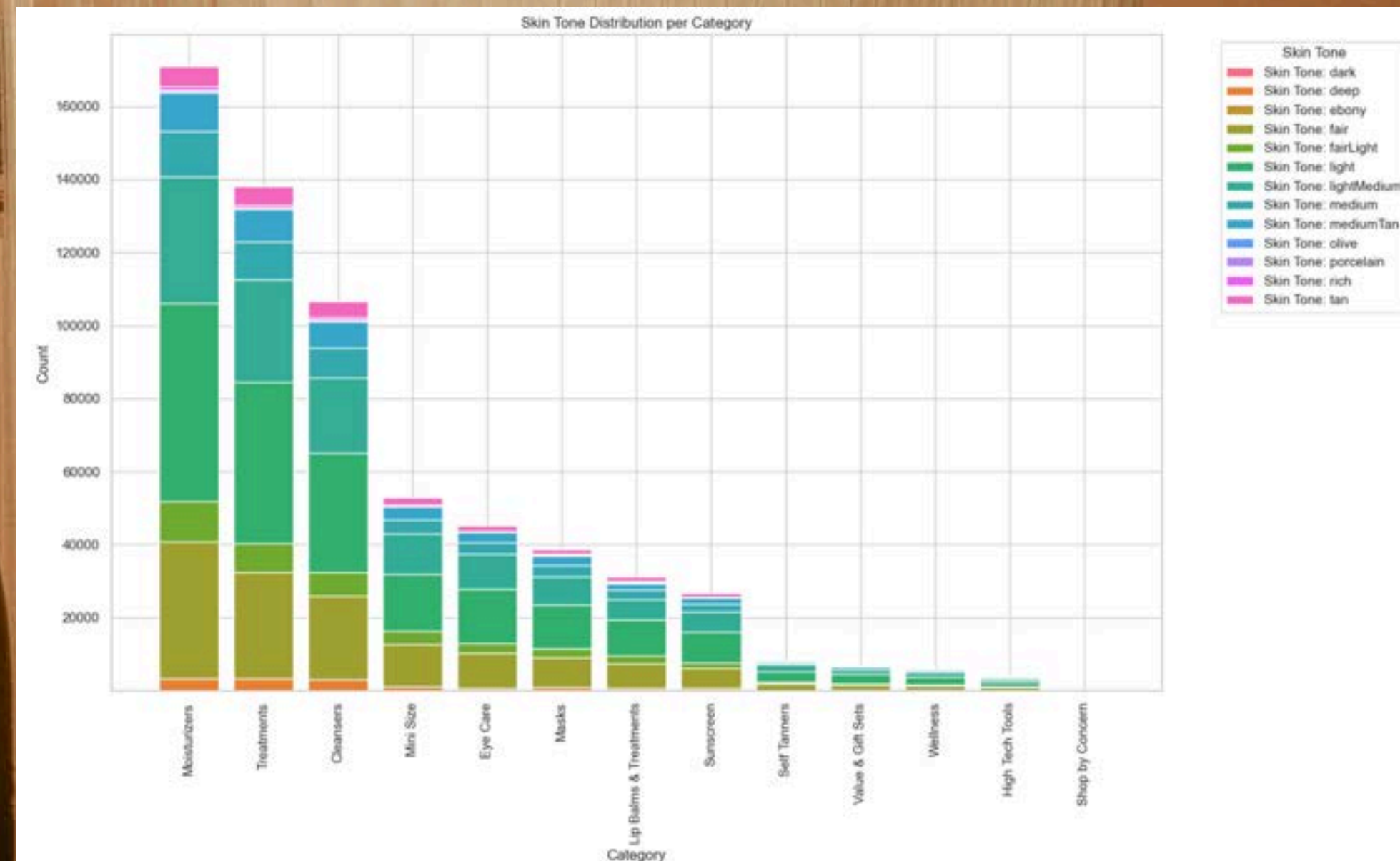Distribution of Recommendation

- Out of 635739 reviews
- 531104 recommended
- 104635 not recommened

# Hypothesis Testing

- H0: There is no significant association between category and skin_tone.
- H1: There is significant association between categroy and skin_tone.

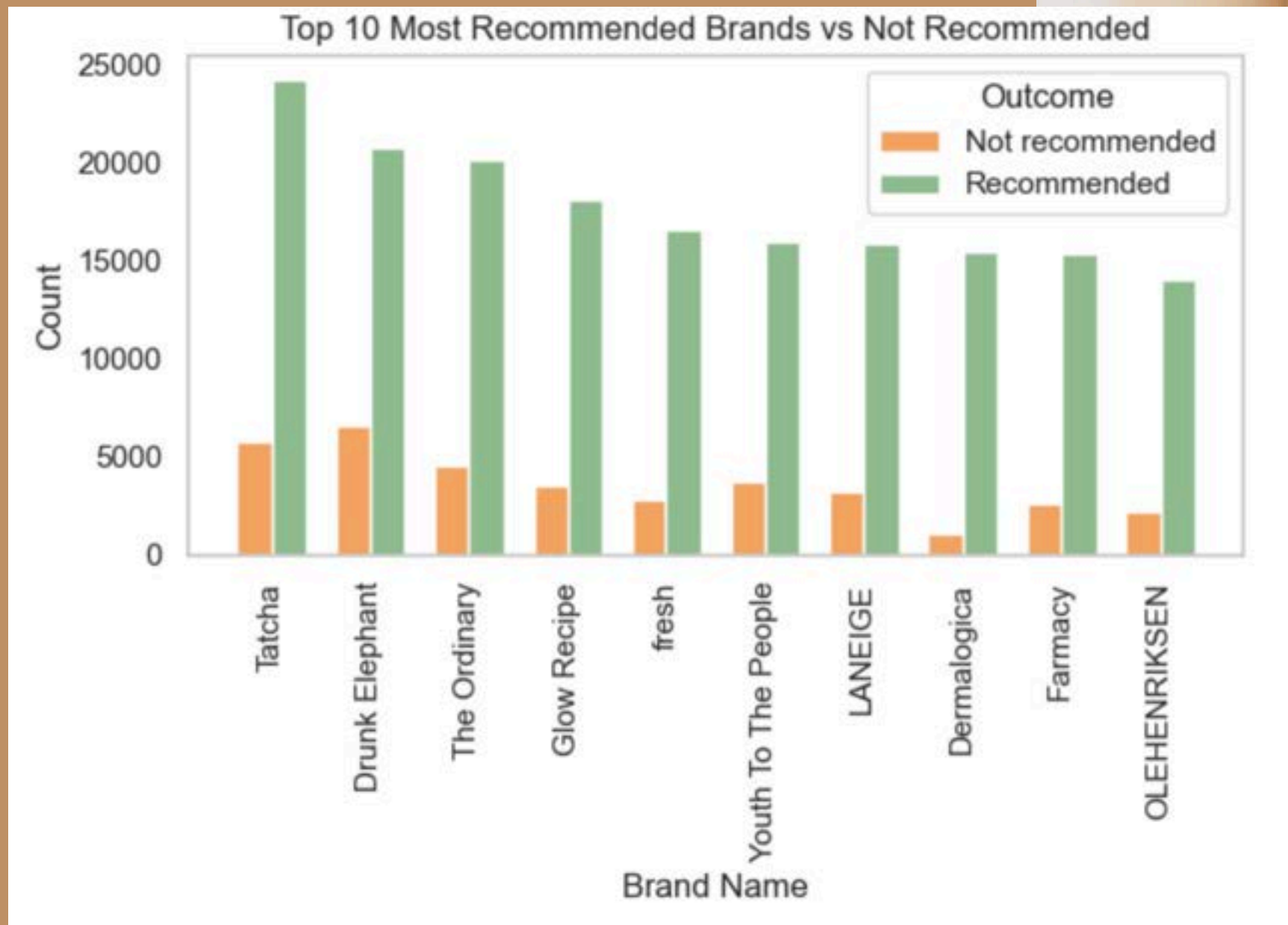- Chi Squared test
- Reject null hypothesis
- Small Cremre's V 0.02


Skin Tone Distribution per Category

# Hypothesis Testing

**SEPHORA**

- H0: There is no significant association between is recommended and brand name.
- H1: There is a significant association between is recommended and brand name.

- Chi - squared test
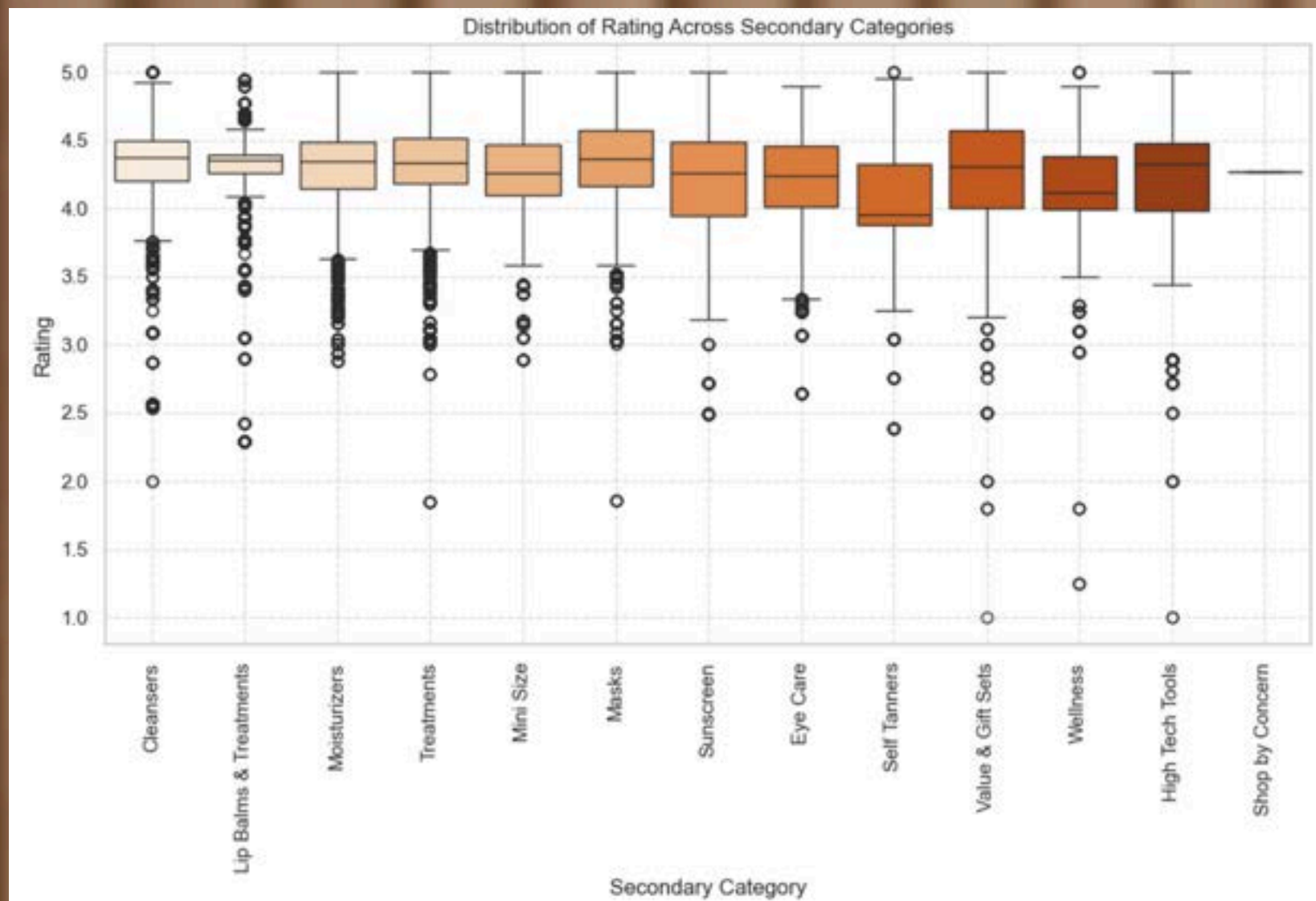- Reject Null Hypothesis
- Small Cremer's V 0.13
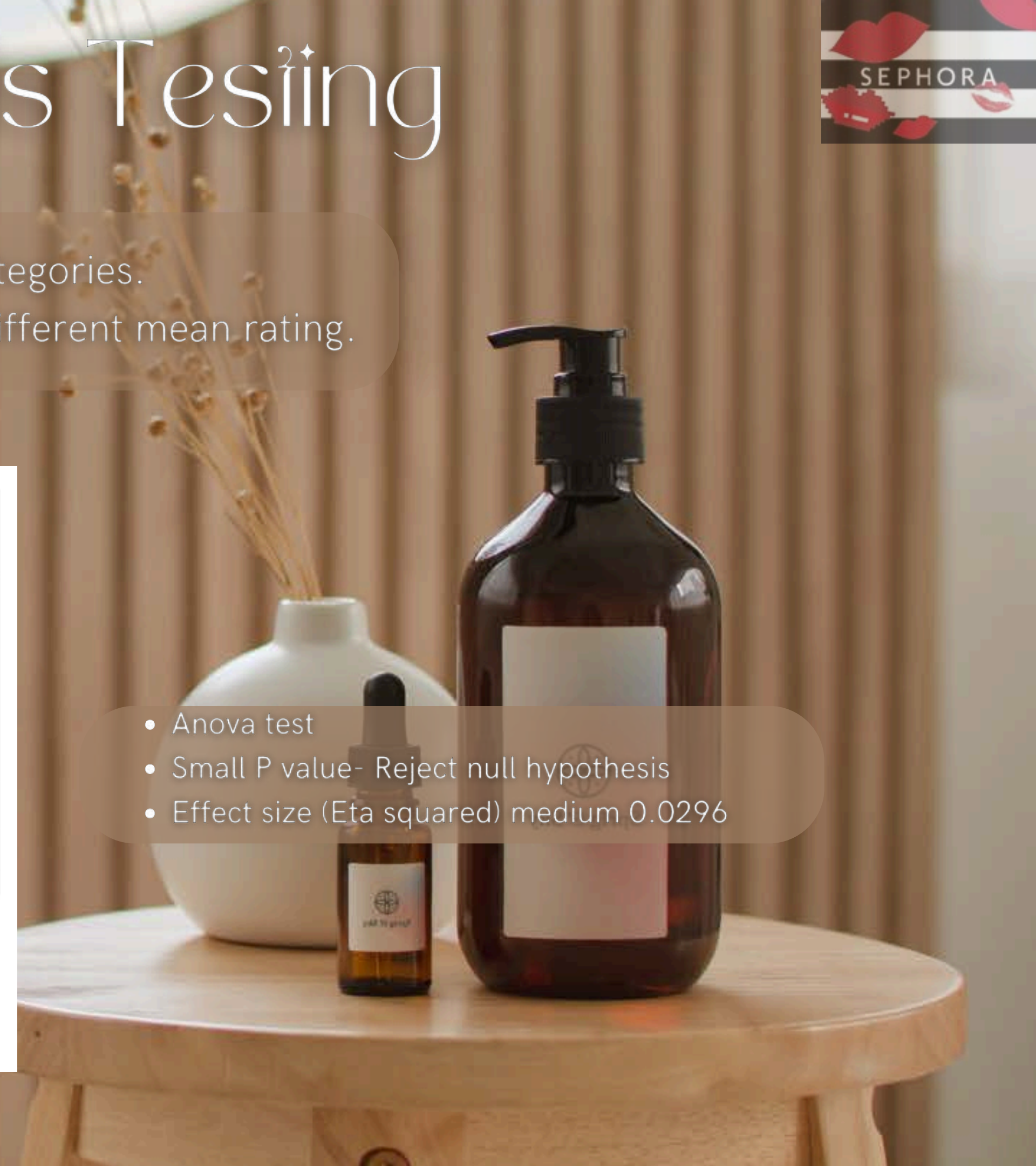


Top 10 Most Recommended Brands vs Not Recommended

# Hypothesis Testing

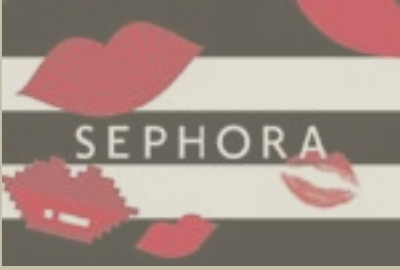- H0: The mean rating is the same across all categories.
- H1: At least one category has a significantly different mean rating.



Distribution of Rating Across Secondary Categories

- Anova test
- Small P value- Reject null hypothesis
- Effect size (Eta squared) medium 0.0296

# Machine Learning Feature selection



Dealing with Multicollinearity



Top Feature Importances

**Selected features:**

rating, skin_tone, total_pos_feedback_count, loves_count, hair_color, skin_type, eye_color, reviews, sales_price, brand_name, category, online_only, review_title, review_text, etc

**Target:**

Is_recommended

# Machine Learning Classification

**Models:**

- Decision Tree
- Logistic Regression
- Random Forest
- Suport Vectore Model
- XGBoost

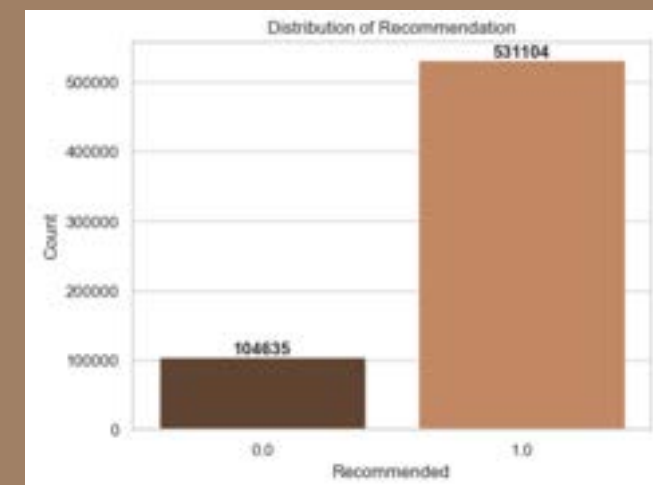| | Model | Test Accuracy | Train Accuracy | Precision (0) | Recall (0) | F1-score (0) | Precision (1) | Recall (1) | F1-score (1) |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Decision Tree | 0.9746 | 0.9757 | 0.89 | 0.96 | 0.92 | 0.99 | 0.98 | 0.98 |
| 1 | Logistic Regression | 0.9446 | 0.9574 | 0.81 | 0.85 | 0.83 | 0.97 | 0.96 | 0.95 |
| 2 | Random Forest | 0.9628 | 0.9643 | 0.92 | 0.82 | 0.87 | 0.97 | 0.99 | 0.98 |
| 3 | SVM | 0.9708 | 0.9725 | 0.88 | 0.95 | 0.91 | 0.99 | 0.97 | 0.97 |
| 4 | XGBoost | 0.9685 | 0.9732 | 0.91 | 0.90 | 0.90 | 0.98 | 0.98 | 0.98 |

# Machine Learning Imbalanced Data

| Model | Precision -0(O) | Recall -0(O) | F1-score -0(O) | Precision -1(O) | Recall -1(O) | F1-score -1(O) | Precision -0(S) | Recall -0(S) | F1-score -0(S) | Precision -1(S) | Recall -1(S) | F1-score -1(S) | Precision -0(T) | Recall -0(T) | F1-score -0(T) | Precision -1(T) | Recall -1 (T) | F1-score -1(T) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Decision Tree | 0.89 | 0.96 | 0.92 | 0.99 | 0.98 | 0.98 | 0.92 | 0.99 | 0.95 | 0.99 | 0.99 | 0.99 | 0.88 | 0.84 | 0.86 | 0.97 | 0.97 | 0.97 |
| Logistic Regression | 0.81 | 0.85 | 0.83 | 0.97 | 0.96 | 0.95 | 0.99 | 0.98 | 0.99 | 0.99 | 0.99 | 0.99 | 0.92 | 0.84 | 0.86 | 0.98 | 0.97 | 0.97 |
| Random Forest | 0.92 | 0.82 | 0.87 | 0.97 | 0.99 | 0.98 | 0.97 | 0.99 | 0.98 | 0.99 | 0.99 | 0.99 | 0.96 | 0.87 | 0.87 | 0.96 | 0.98 | 0.98 |
| SVM | 0.88 | 0.95 | 0.91 | 0.99 | 0.97 | 0.97 | 0.98 | 0.97 | 0.97 | 0.99 | 0.99 | 0.99 | 0.94 | 0.92 | 0.93 | 0.98 | 0.98 | 0.98 |
| XGBoost | 0.91 | 0.90 | 0.90 | 0.98 | 0.98 | 0.98 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.94 | 0.93 | 0.93 | 0.99 | 0.99 | 0.99 |

**Models:**

- Decision Tree
- Logistic Regression
- Random Forest: Good performance without Resampling
- SVM: Best Performer with Tomek
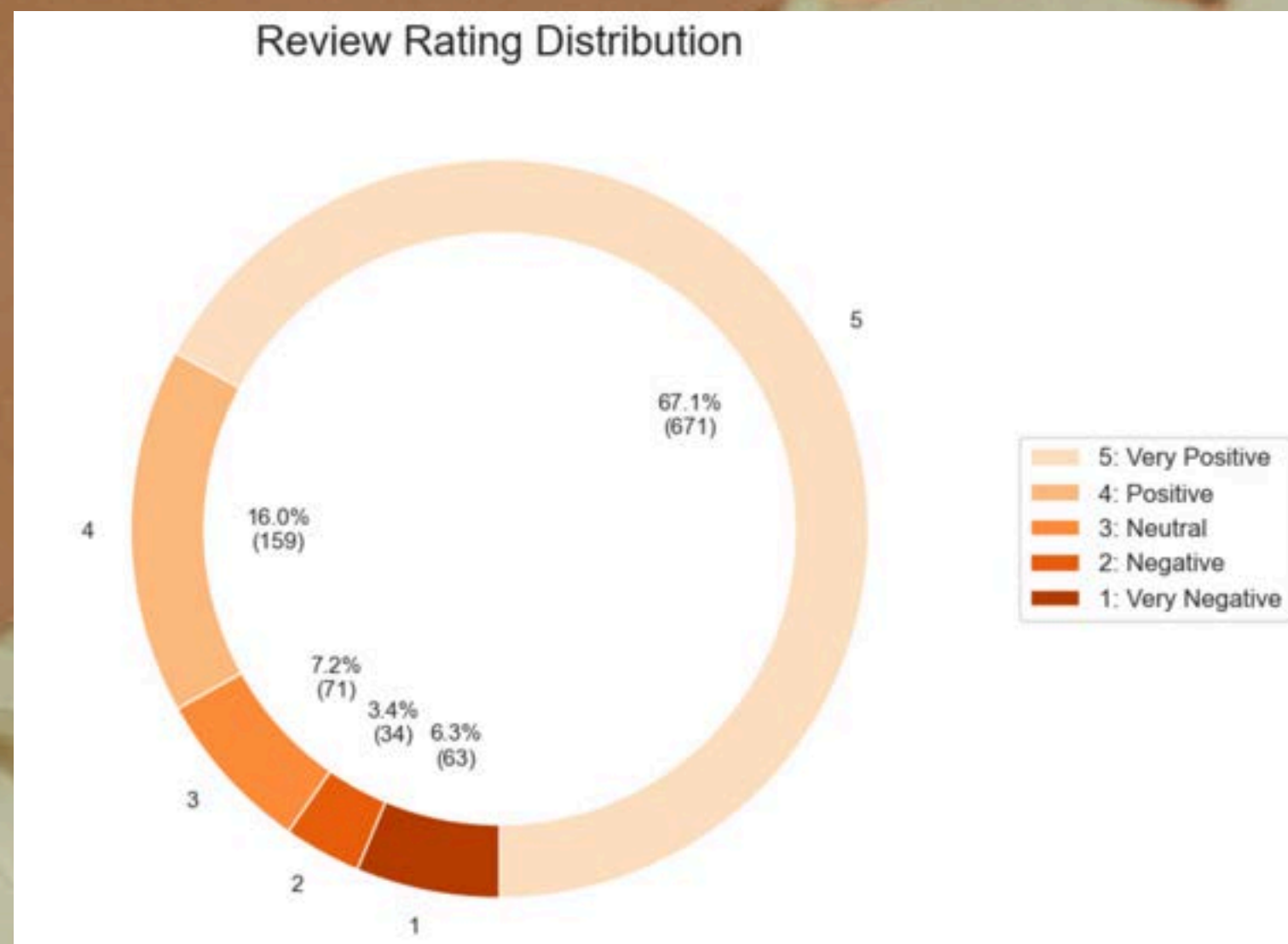- XGBoost: Best Performer overll with SMOTE
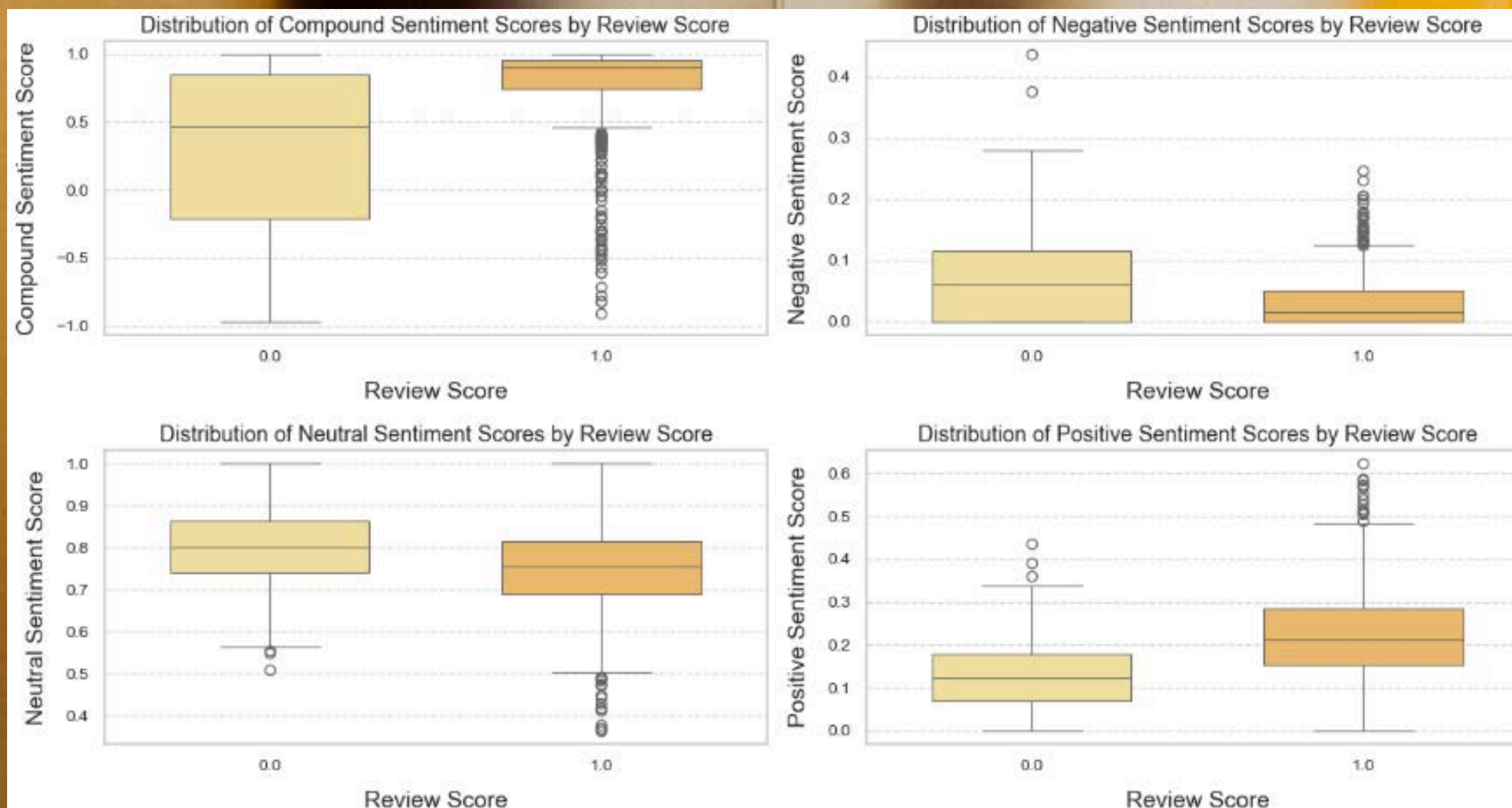
# Review Rating Disirbuiion



- Majority are positive 83%
- 9.7% Negative
- Smallest fraction of 7.2% neutrals

Review Rating Distribution

67.1%
(671)

16.0%
(159)

7.2%
(71)

3.4%
(34)

6.3%
(63)

5

4

3

2

1

5: Very Positive
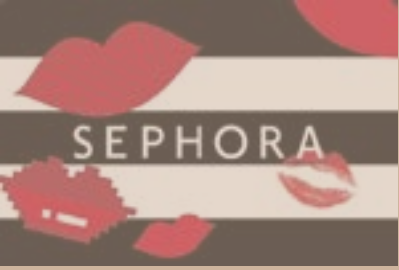4: Positive
3: Neutral
2: Negative
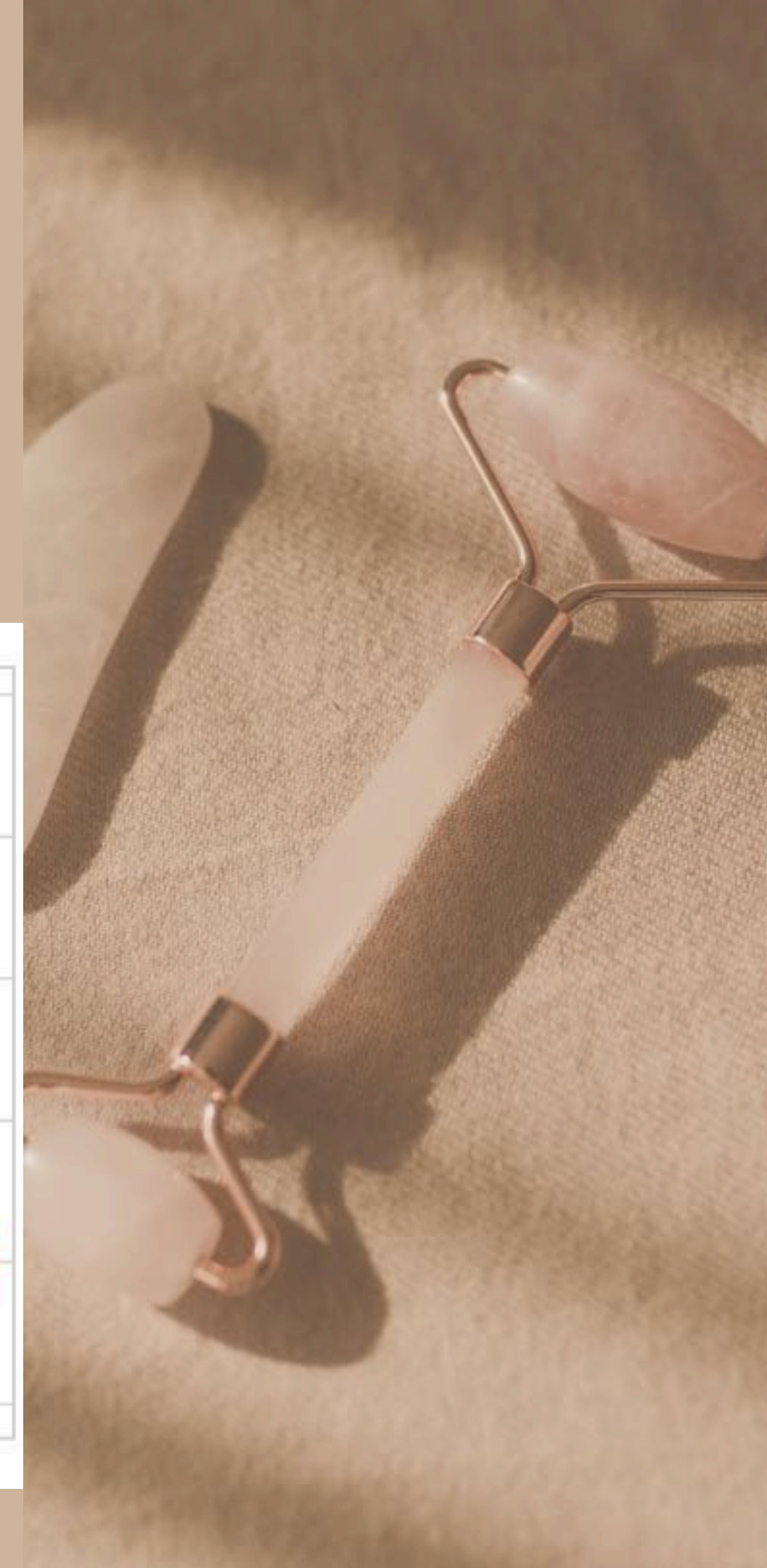1: Very Negative

# Distribution of sentiment scores by is_recommended
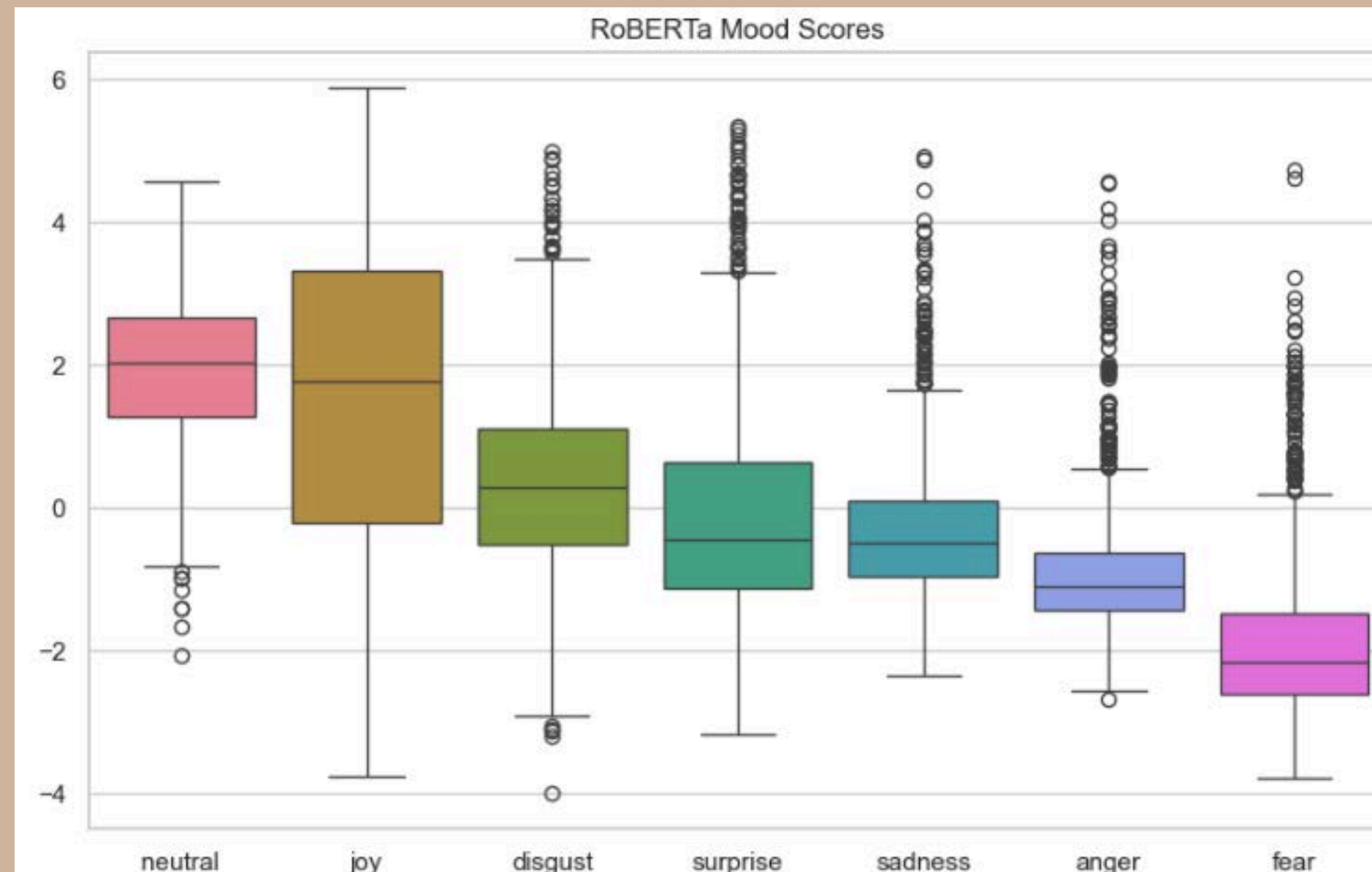
- Overall higher score for recommendation (Compound & Positive)
- Higher score in Non_recommended (Negative & Neutral)



Distribution of Compound Sentiment Scores by Review Score
Distribution of Negative Sentiment Scores by Review Score
Distribution of Neutral Sentiment Scores by Review Score
Distribution of Positive Sentiment Scores by Review Score

# Sentiment Analysis RoBERTa

- Identifying specific consumers' emotion and mood detection
- Neutral and Joy highest median illustrating customer satisfaction
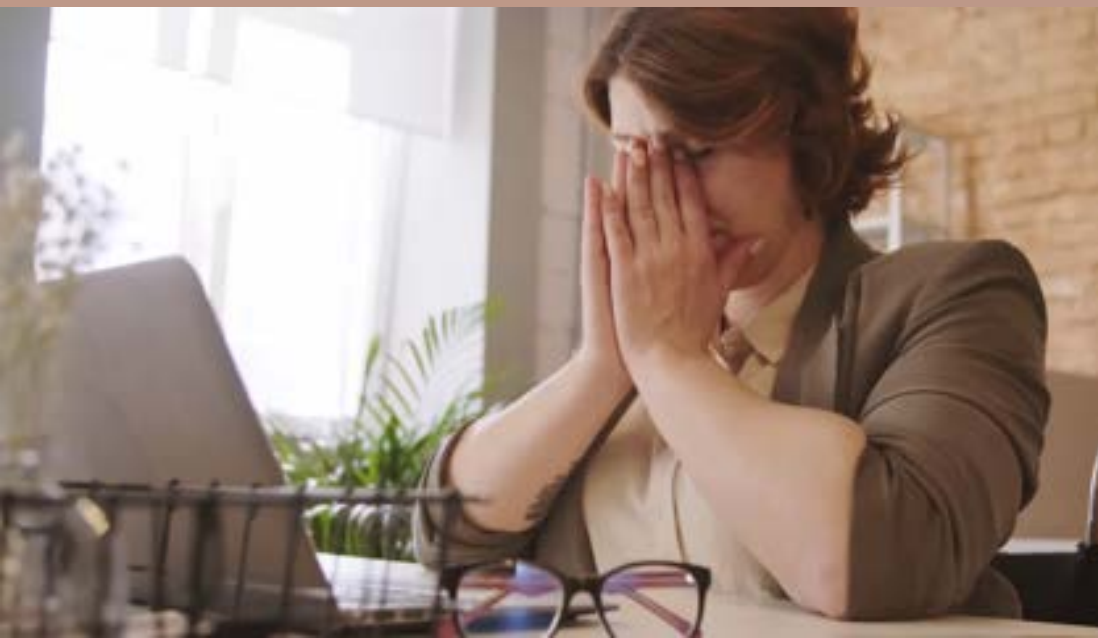
# Challenges & Solutions

## Challenges

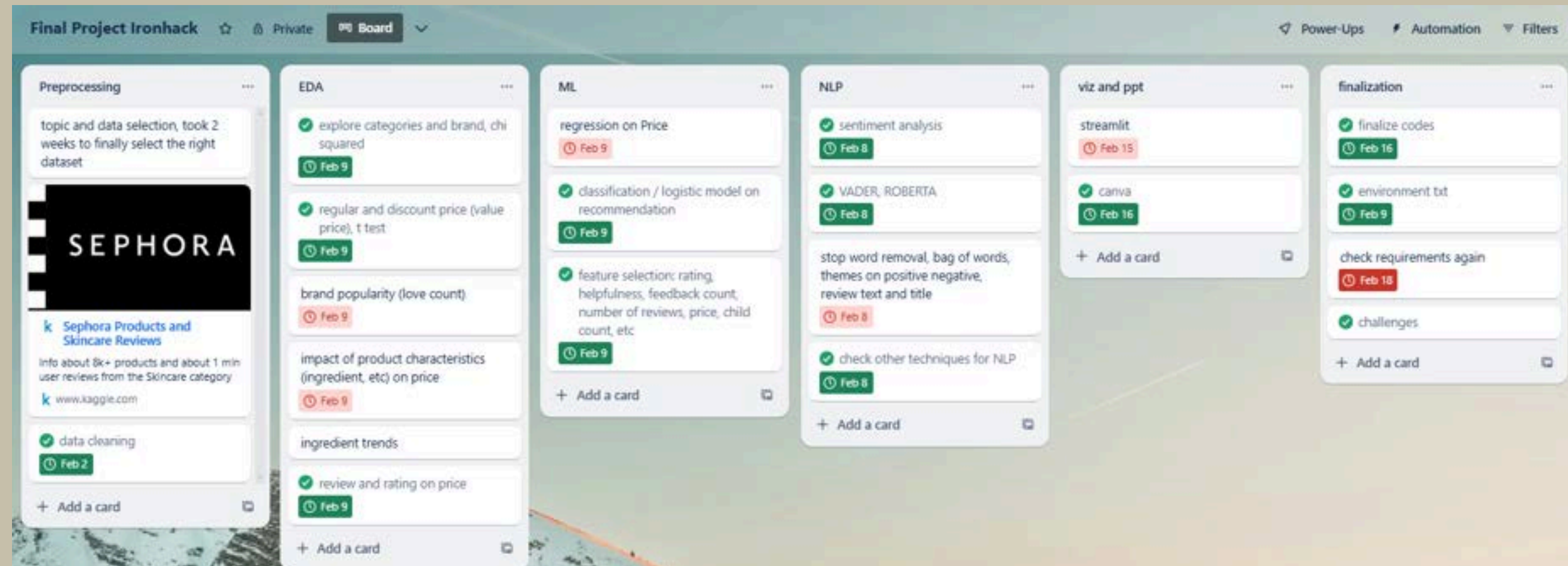- Large datasize long processing time
- Memory error
- High-dimensionality

## Solution

- Smaller sample size
- Textblob for encoding

# Planning

# Thank you