# Predicting the 'Usefulness' of Yelp Reviews

**Presented by Gyujin Seo | UCLA | Major in Statistics and Data Science**

# Table of Contents

# 1. Introduction

▶ **Background**

   ▷ Today, we live an era where users' experiences are highly valued and expressing them is crucial.

   ▷ Companies recognize the importance of such experiential data(review data) and utilize them as significant marketing resources.

   ▷ Through the "yelp" website, which provides consumer review datasets, these data are collected.

▶ **Problem**

   ▷ Excessive time spent on purchase decision-making due to the ambiguity of consumer review data reliability.

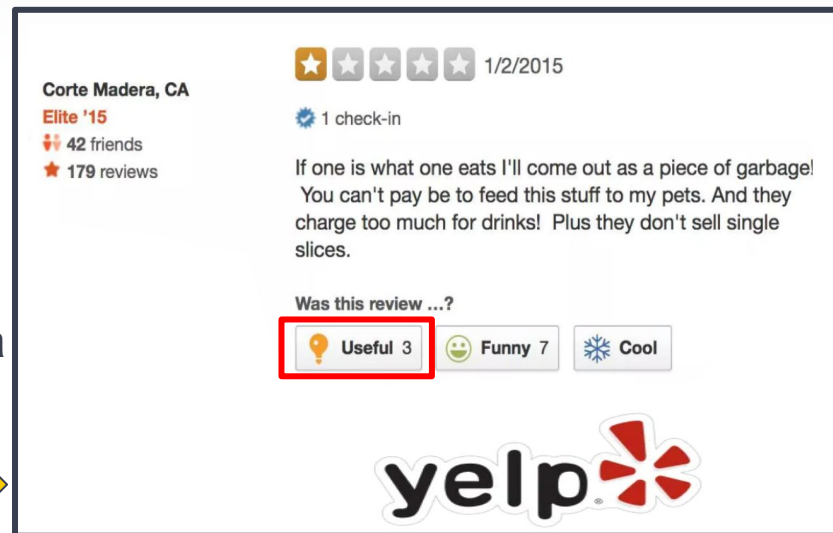   ▷ Companies spending excessive time determining the usefulness of review data.

# 1. Introduction(Cont'd)

▶ **Aim of my Project**

▷ To predict the usefulness of a review prior to user interaction.

▶ **Yelp Overview**

▷ Yelp is a consumer-driven platform providing business information, photos, and user reviews.

▷ The **'useful'** button on Yelp, which plays a critical role in guiding user decisions.

# 1. Introduction(Cont'd)

- **Programming language**
  - R

- **Data overview**
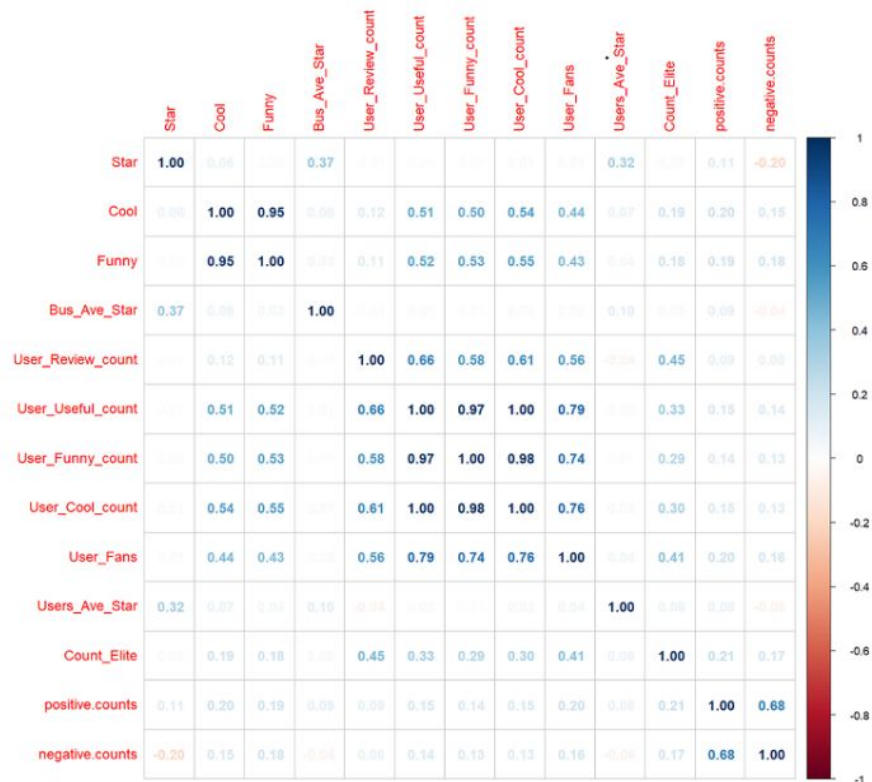  - It has 53,845 unique observations, with 17 variables.
  - For the variables, I have User ID, Business ID, Star, Useful, Cool, Funny, Review, State, City, Business average count, User review count, User useful count, Userfunny count, User cool count, Elite, User fans, and User average star.

# 1. Introduction(Cont'd)

| Variables | Description |
|---|---|
| User_id | A random code with letters, numbers and special characters "_" and "" that identifies each user |
| Bus_id | A random code with letters, numbers and special characters "_" and "" that identifies the business |
| Star | In a scale of 1 to 5, the star represents different levels of satisfaction a user may have for the service at the restaurant, where each number represents the following: 1 - Not good, 2 - Could've been better, 3 - OK, 4 - Good, 5 - Great |
| Useful | The number of "Useful" votes a review gets from other users |
| Cool | The number of "Cool" votes a review gets from other users |
| Funny | The number of "Funny" votes a review gets |
| Review | A detailed description of the user's experience at the restaurant |
| State | The business state |
| City | The business city |
| Bus_Ave_Star | The averaged stars of a business |
| User_Review_count | The total number of reviews by this user |
| User_Useful_count | The total number of "Useful" votes of all the reviews by this user |
| User_Funny_count | The total number of "Funny" votes of all the reviews by this user |
| User_Cool_count | The total number of "Cool" votes of all the reviews by this user |
| Elite | The years a user is selected to be an Elite |
| User_Fan | The number of followers a user has |
| Users_Ave_Star | The average of all the ratings a user gives |

# 2. Pre-processing

- Transforming Useful to binary: "Not Useful" (0) if less than or equal to 1, "Useful" (1) otherwise.

- Review analysis with NLP: Counted positive and negative words in each review.

- Eliminating highly correlated variables: Dropped Funny, User Funny count, User Cool count due to high correlation (>90%).
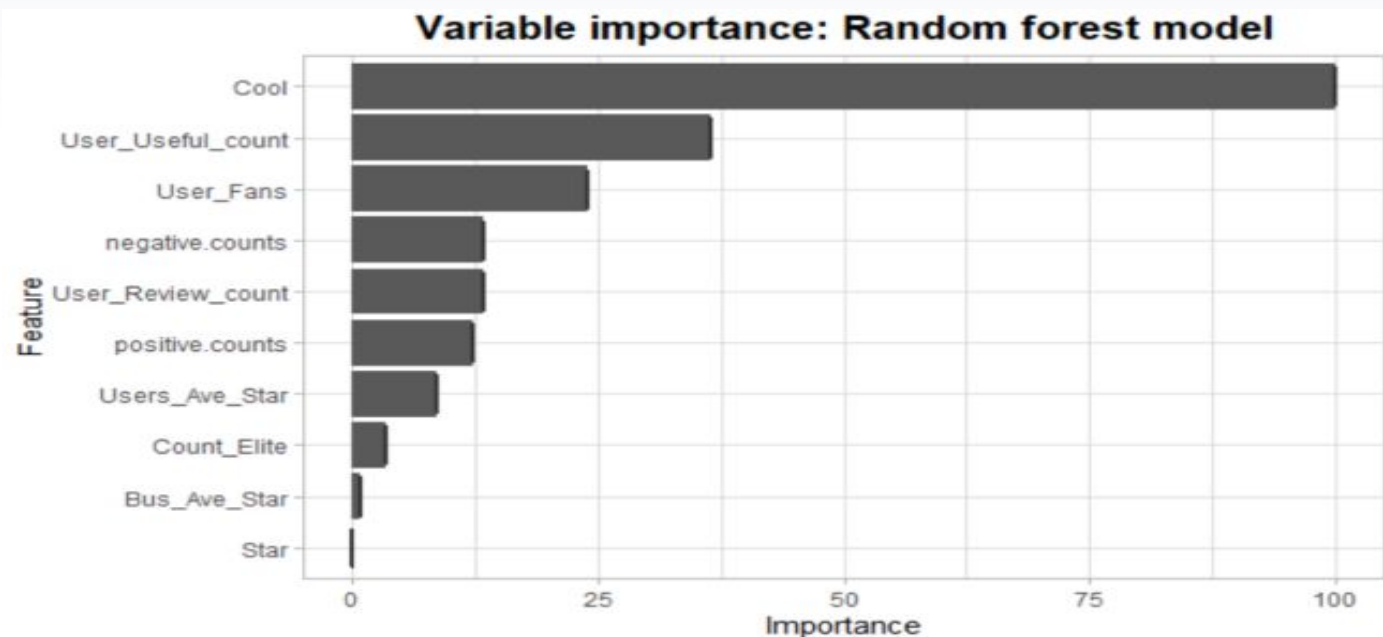
# 3. Modeling

## 3.1 Model Validation

- **Data Partitioning: 70% Training Data, 30% Test Data.**
- **Cross-Validation: Applied 10-fold Cross-Validation for more accurate performance estimation.**

## 3.2 Methods

- **5 Machine Learning Methods Applied:**
  - **Random Forests**
  - **Logistic Regression**
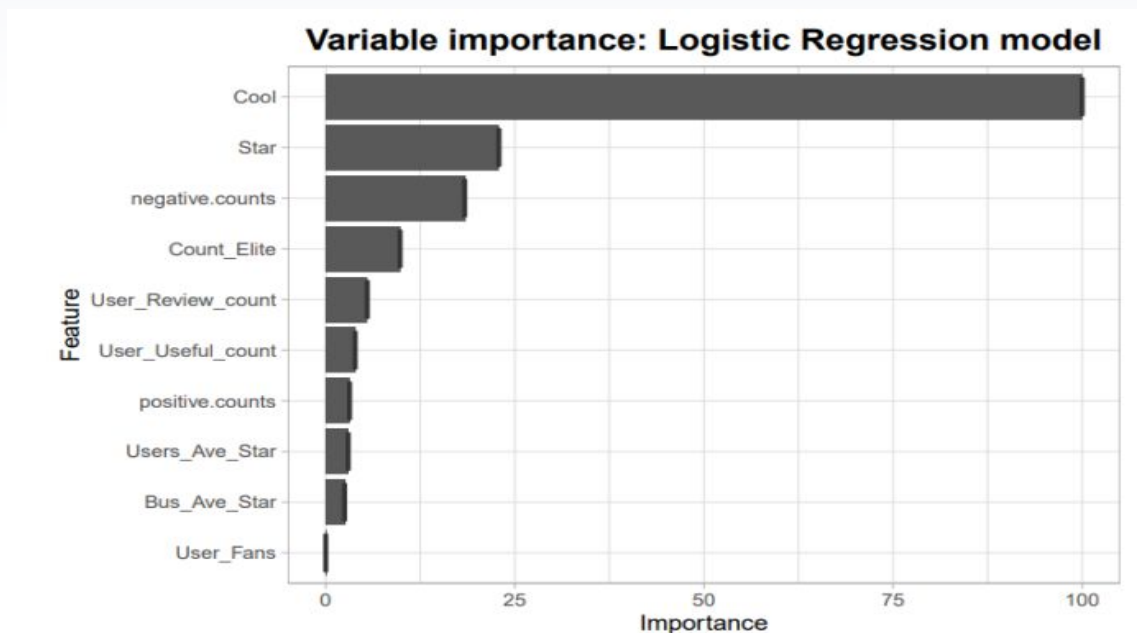  - **LASSO Regression**
  - **Decision Tree**
  - **K-NN**
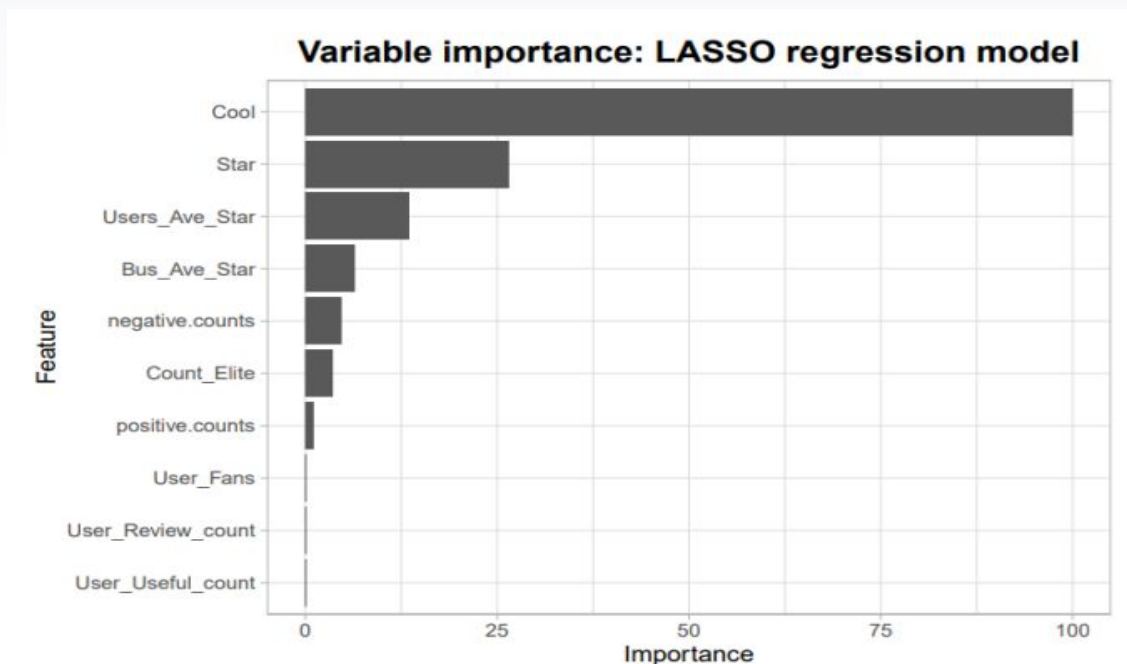
# 3. Modeling(Cont'd)

## 3.2 Methods - Random Forest



Variable importance: Random forest model

# 3. Modeling(Cont'd)

## 3.2 Methods - Logistic Regression



Variable importance: Logistic Regression model

# 3. Modeling(Cont'd)

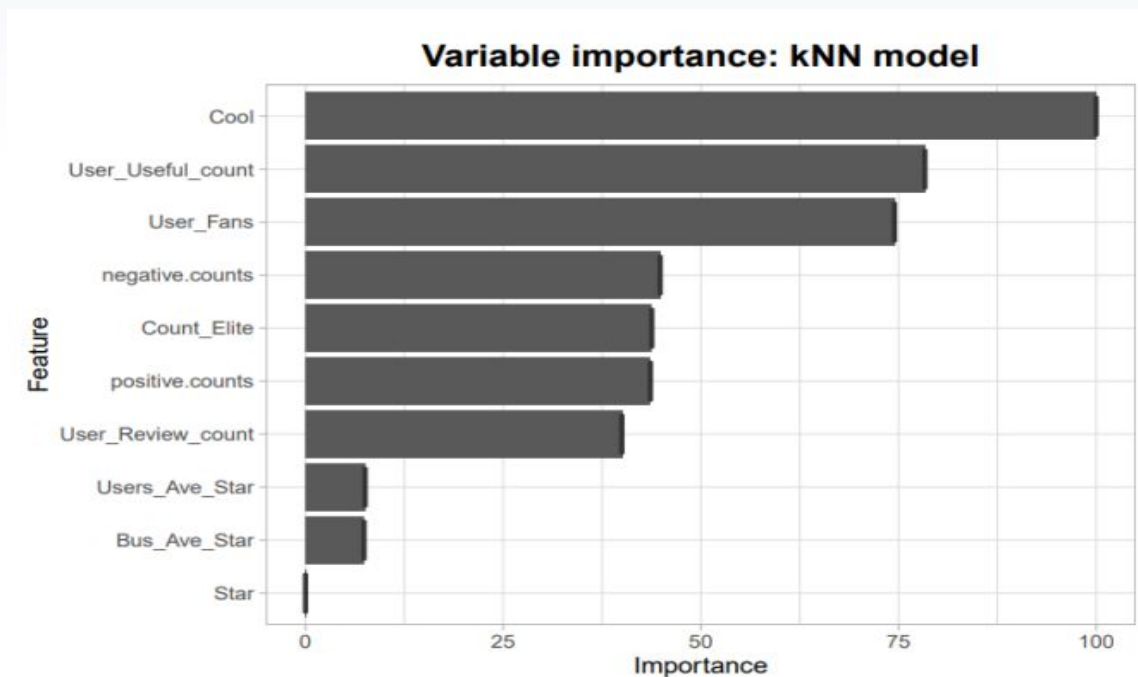## 3.2 Methods - LASSO Regression



Variable importance: LASSO regression model

# 3. Modeling(Cont'd)

## 3.2 Methods - Decision Tree

# 3. Modeling(Cont'd)

## 3.2 Methods - K-NN

# 4. Results and Conclusion



ROC curves

| Method | AUC | Accuracy | F1_score |
|---|---|---|---|
| Random Forest | 0.9158 | 0.8759 | 0.7563 |
| Logistic Regression | 0.9096 | 0.8760 | 0.7565 |
| LASSO Regression | 0.9099 | 0.8769 | 0.7564 |
| Decision Tree | 0.8456 | 0.8723 | 0.7463 |
| k-NN | 0.8460 | 0.8324 | 0.6560 |

▶ **'Random Forest' has an AUC score of 91.58% that shows the model is learning the data well enough.**

# 4. Results and Conclusion(cont'd)

- **Consumer Perspective**
  - By predicting the usefulness of reviews, we can secure the trust of buyers and expedite their purchase decision-making process.

- **Restaurant Perspective**
  - Predicting the 'usefulness' of reviews not only allows for swift utilization as marketing material, but can also enhance customer service satisfaction.

# 5. Future Research

▶ **Research to derive meaningful patterns and insights from the given dataset using various data mining techniques.**

  ▷ Using algorithms like K-means clustering for customer segmentation.

▶ **Research on the development of predictive models using machine learning algorithms.**

  ▷ Used for trend prediction, user behavior prediction, recommendation systems, etc.

▶ **Research to improve model performance through optimization theory and techniques.**

  ▷ Solving issues such as model complexity, overfitting, and underfitting.

# Thank you!