# Amazon Fake Review Detection using various ML Model
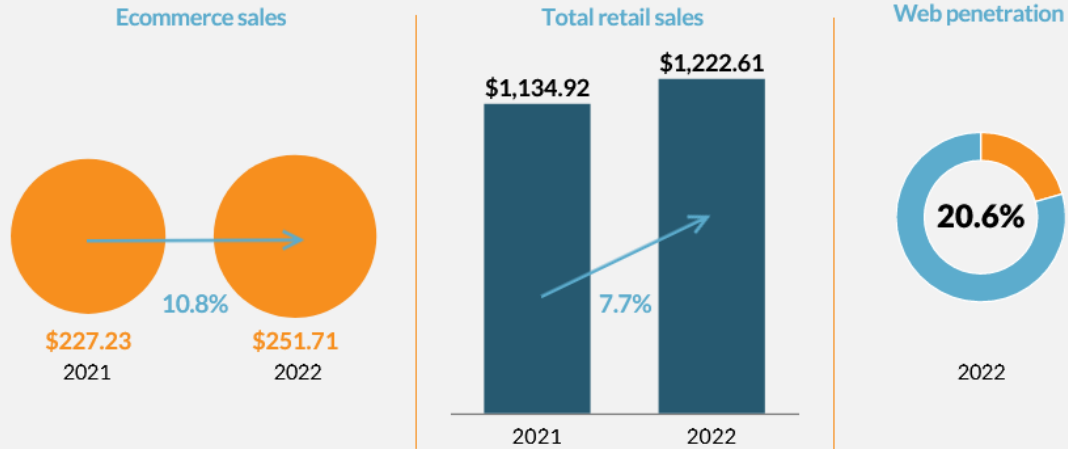
DONGJUN CHO

605.744

# Ecommerce Growth
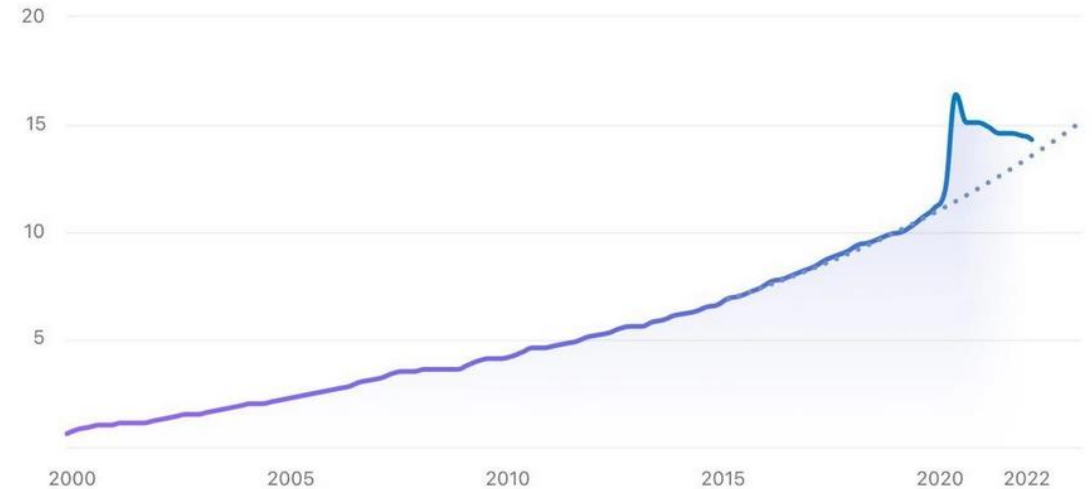


US retail landscape during Q3 2022, in $billions

Ecommerce sales: $227.23 (2021) → $251.71 (2022), 10.8%

Total retail sales: $1,134.92 (2021) → $1,222.61 (2022), 7.7%

Web penetration: 20.6% (2022)



US ecommerce adoption growth rate

% OF ADDRESSABLE RETAIL

U.S. ecommerce grew 10.8% in Q3 2022.

- E-commerce sales accounted for around 14% of retail sales

# About Review

Reviews are key to the decision-making process, helping customers to get a better idea about the product

95% of consumers read online reviews before they shop
- ◦ 58% say they would pay more for the products of a brand with good reviews.

# Related Works

In 2008, Jindal and Liu identified spam by analyzing the number of feedbacks, length of review titles, sorting order by date, positives used repeatedly in content, percentage of negative words, and review evaluation ratings for amazon reviews.

In 2021, Alsubari used Chicago hotel reviews (1600 hotel reviews, 800 fake, 800 real review), and predicted fake reviews with an accuracy of 93% and 95% with a review classifier using machine learning (SVM, Random Forest).

# Fake Review

These reviews produced by persons who have not personally experienced the subjects of the reviews are called spam reviews; spam reviews might also be called fake reviews, non-genuine reviews, or fraudulent reviews.

◦ Detection of review spam: A survey

# Project Idea

Using supervised machine learning technique to train model to detect fake review based on labeled dataset.
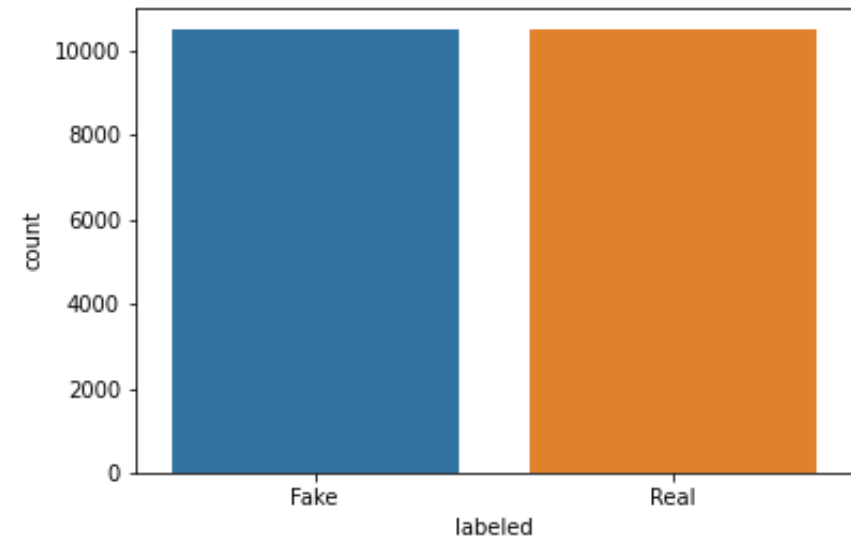
# Dataset

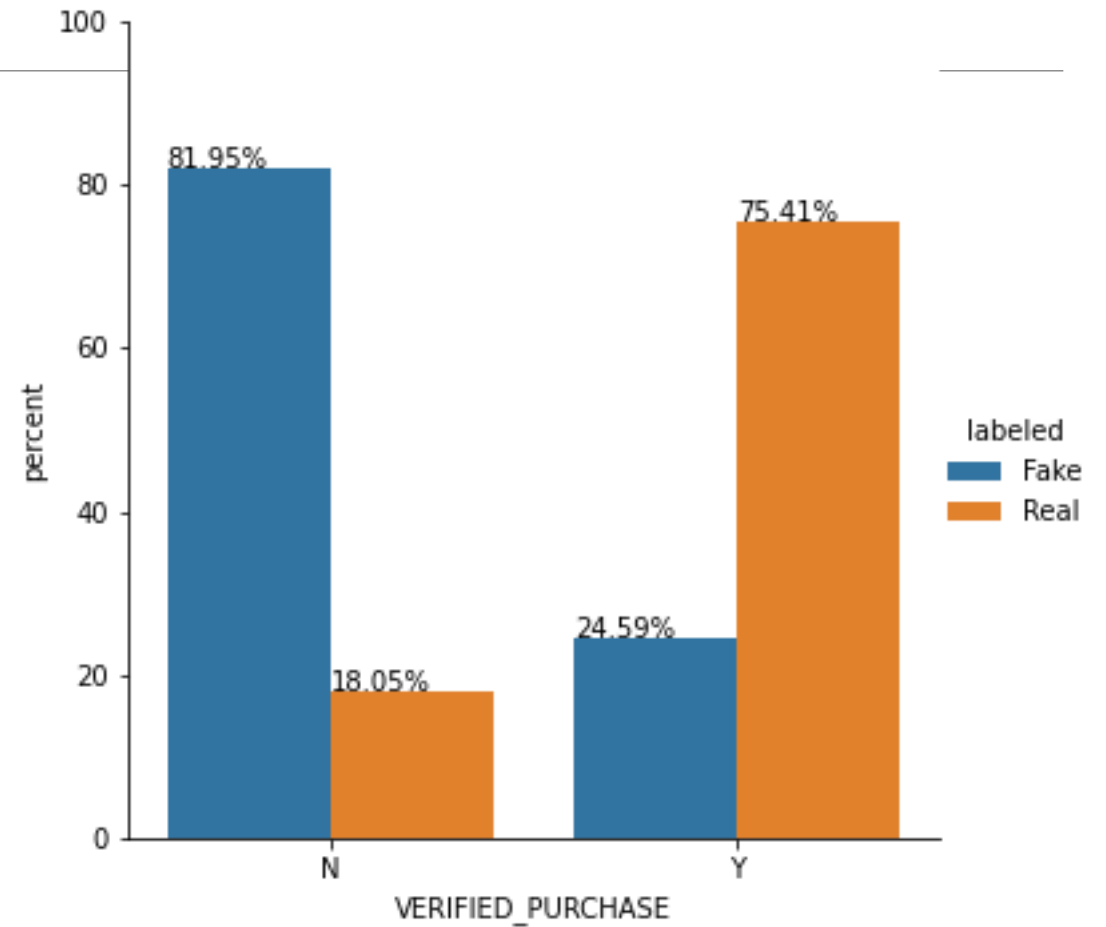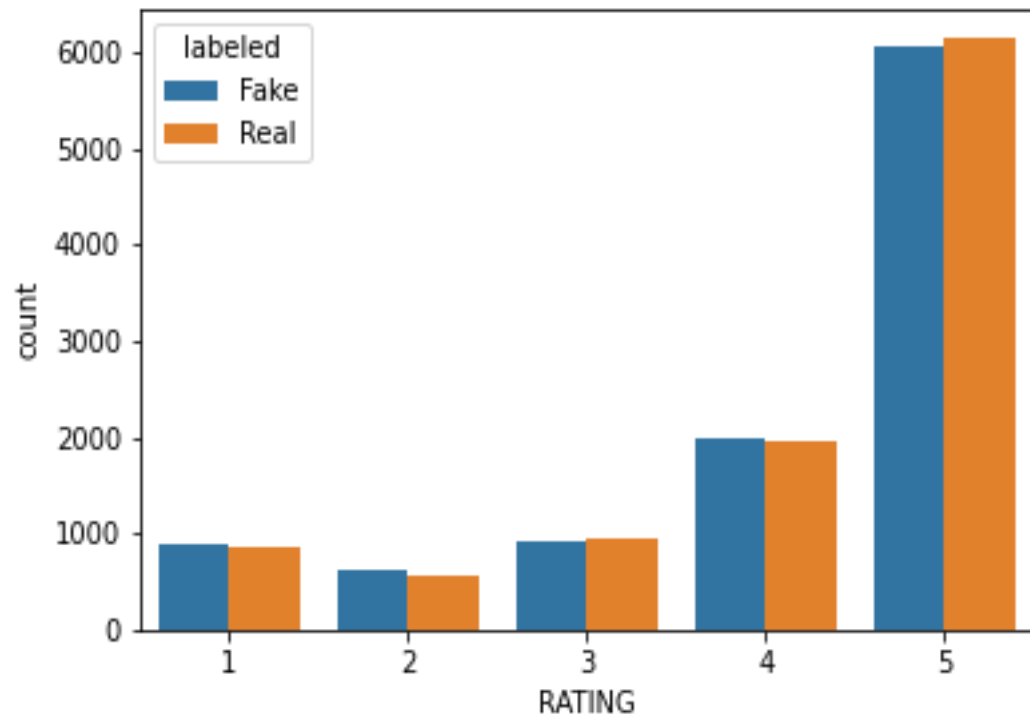Used the dataset from Kaggle and previous research

Dataset of 20000 reviews
◦ 10000 fake reviews
◦ 10000 Real reviews

Reviews about amazon product reviews.

# Data Analysis
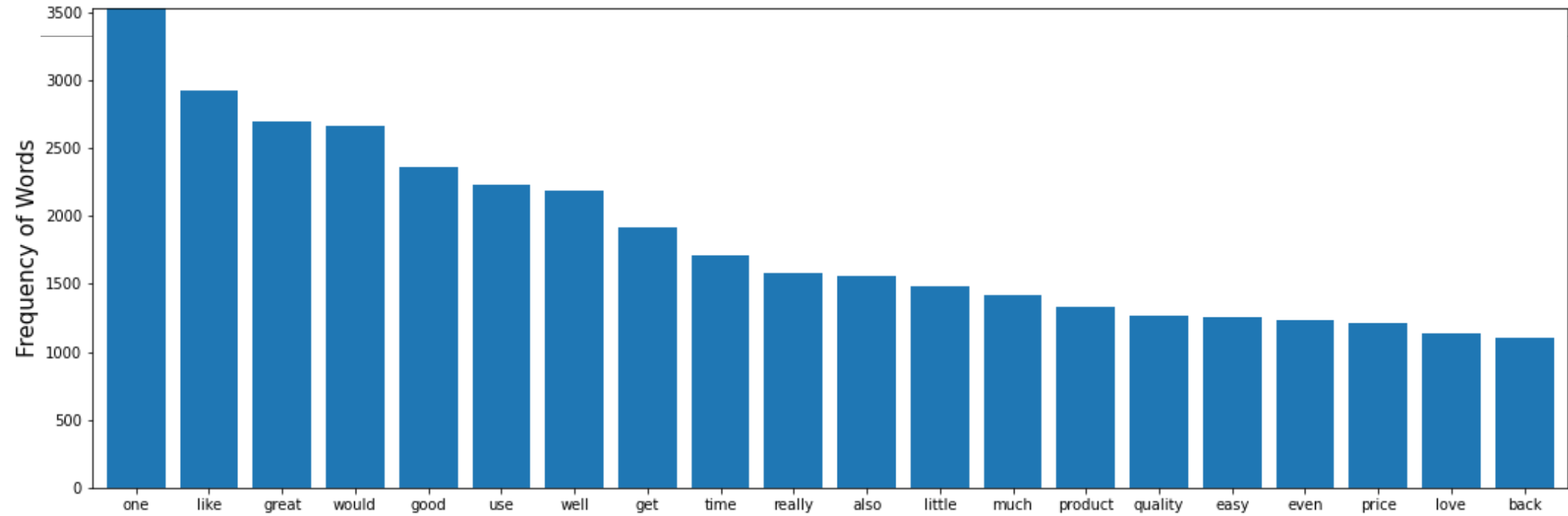
# Data Analysis

Average Review Title Length
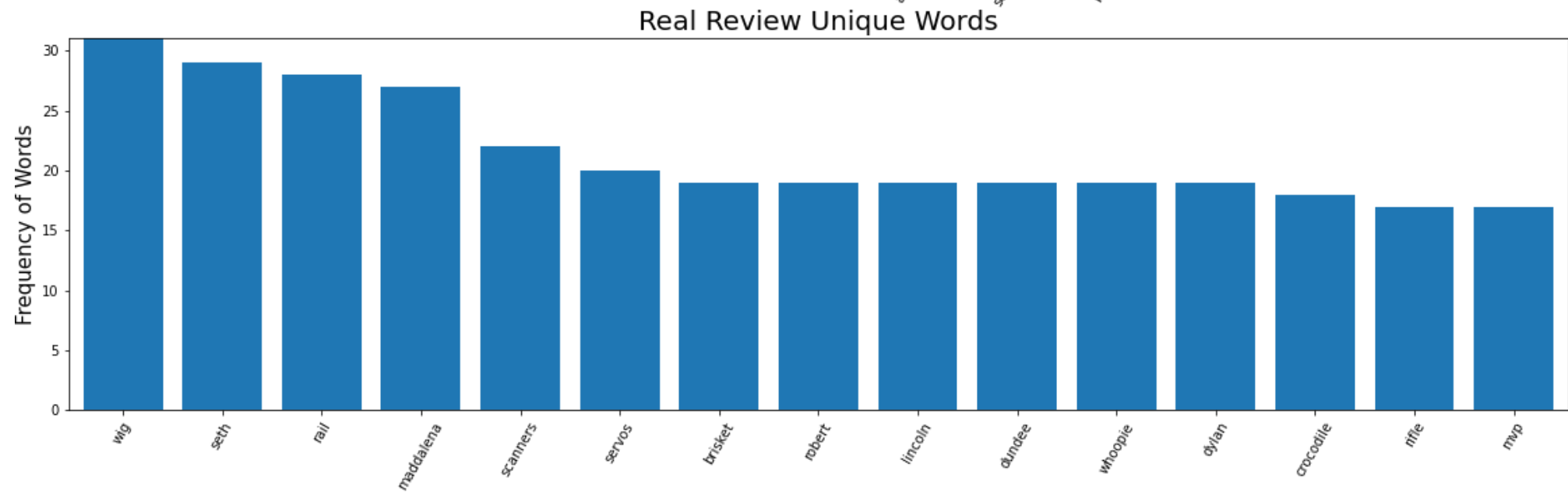- Fake Review: 2.73
- Real Review: 3.18

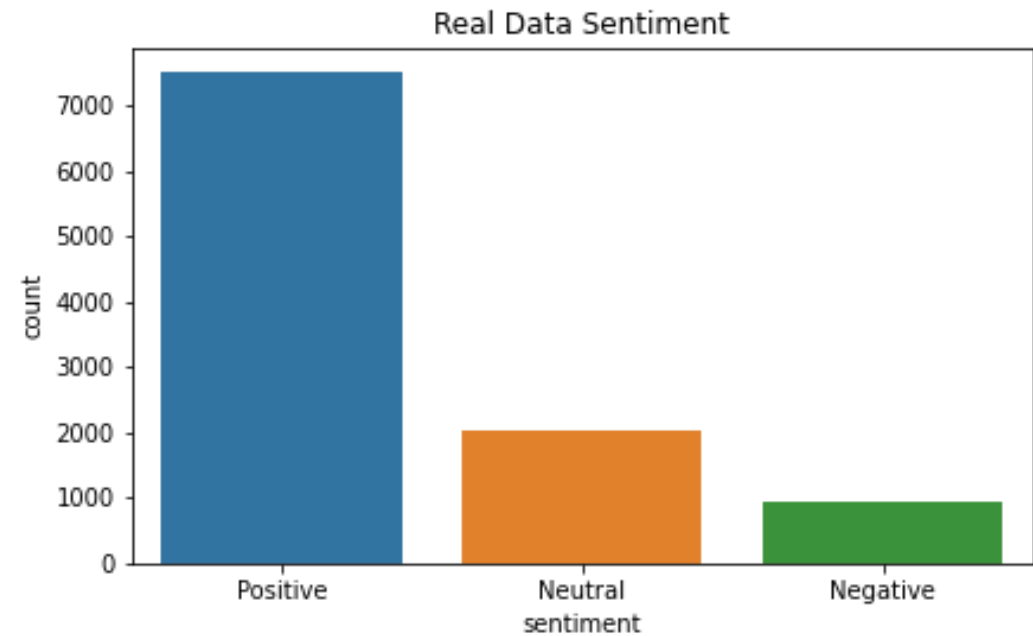Average Review Length
- Fake Review: 28.53
- Real Review: 38.76

Common words that occurs in both real and fake text

Fake Review Unique Words

Real Review Unique Words

# Sentiment Analysis

**Subjectivity**

0 ——————— +1

Objective (fact)    Subjective (opinion)

**Polarity**

-1 ——————— +1

Negative    Positive

# Steps

Dataset

Data Processing

Feature Extraction (TF-IDF)

Data Splitting

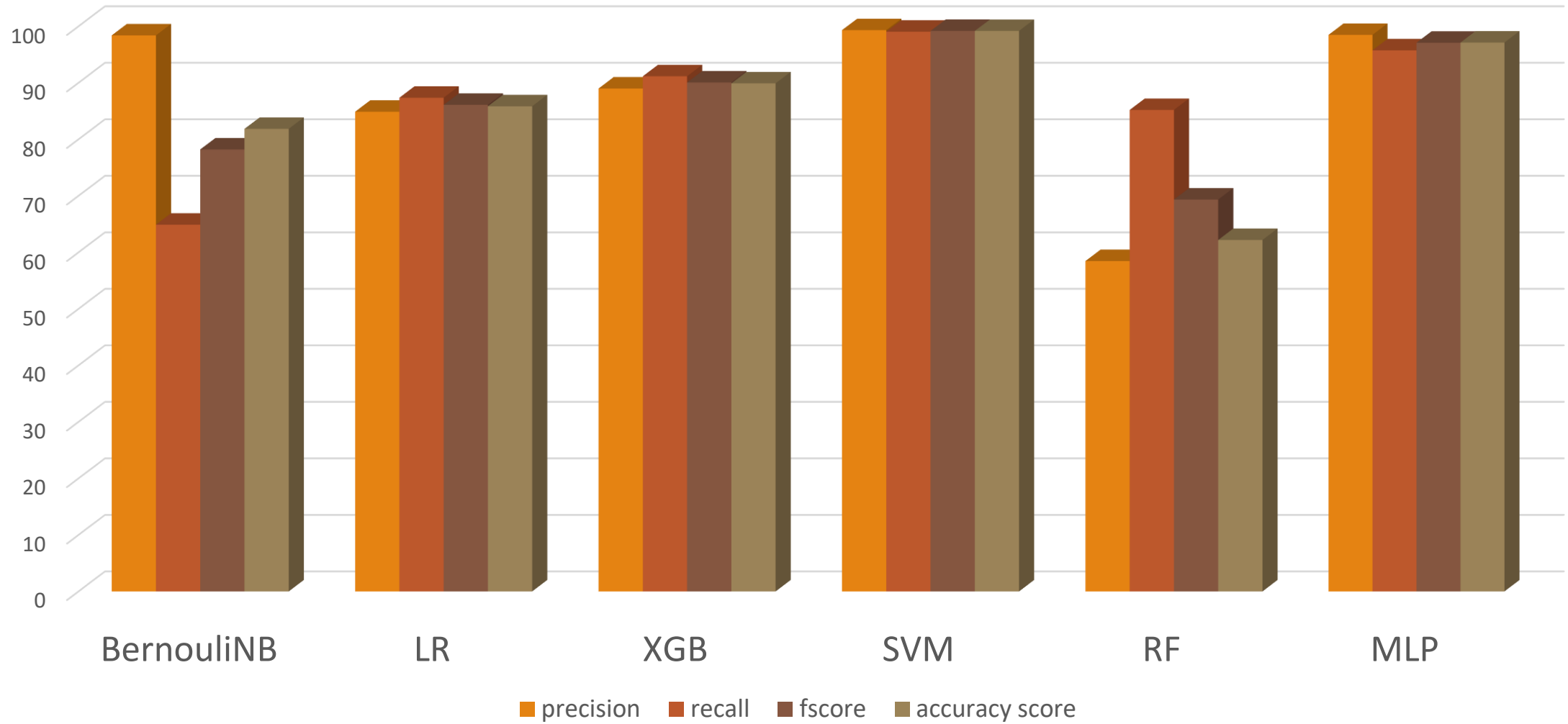Supervised Machine Learning Techniques

Performance Evaluation

# Machine Learning Models

1. Naïve Bayes Classifier

2. Logistic Regression

3. XGBoost Classifier

4. Support Vector Machine

5. Random Forest Classifier

6. Multi-layer Perceptron classifier (MLP)

# Experiment Result

|  |  | Precision (%) | Recall (%) | F1 Score (%) | Accuracy Score (%) |
|---|---|---|---|---|---|
| BernouliNB | Train | 90.281 | 54.329 | 67.836 | 74.273 |
|  | Test | 98.389 | 64.896 | 78.207 | 81.848 |
| LR | Train | 78.791 | 81.004 | 79.882 | 79.625 |
|  | Test | 84.882 | 87.362 | 86.104 | 85.848 |
| XGB | Train | 81.986 | 82.289 | 82.137 | 82.127 |
|  | Test | 88.996 | 91.157 | 90.064 | 89.905 |
| SVM | Train | 94.989 | 94.24 | 94.613 | 94.641 |
|  | Test | 99.315 | 99.051 | 99.183 | 99.181 |
| RF | Train | 65.66 | 52.537 | 58.37 | 62.578 |
|  | Test | 58.464 | 85.199 | 69.344 | 62.19 |
| MLP | Train | 96.973 | 96.529 | 96.75 | 96.762 |
|  | Test | 98.478 | 95.75 | 97.094 | 97.124 |

ML Classification Results

# Conclusion

SVM
- ◦ Text Categorization with Support Vector Machines: Learning with Many Relevant Features - Thorsten Joachims

Future Plans / Works
- ◦ Larger dataset
- ◦ Testing new models

# Reference

1. Heydari, A., Tavakoli, M. A., Salim, N., & Heydari, Z. (2015). Detection of review spam: A survey. Expert Systems with Applications, 42(7), 3634-3642.

2. Jindal, N., & Liu, B. (2008). Opinion spam and analysis. Proceedings of the 2008 International Conference on Web Search and Data Mining, Palo Alto, California, USA, ACM, 219-230.

3. Alsubari, Saleh & Deshmukh, Sachin & Alqarni, Ahmed & Alsharif, Nizar & Aldhyani, Theyazn & Fawaz, Waselallah & Alsaade, & Khalaf, Osamah. (2021). Data Analytics for the Identification of Fake Reviews Using Supervised Learning. Computers, Materials and Continua. 70. 10.32604/cmc.2022.019625.

4. Joachims, Thorsten. "Text Categorization with Support Vector Machines: Learning with Many Relevant Features." European Conference on Machine Learning (1998).

# Thanks