

A

Course End Project Report on
Instagram Data Analysis: A Comprehensive Study
of User Engagement and Content Performance
Is submitted in partial fulfillment of the Requirements for the Award of CIE of
DATA ANALYSIS AND VISUALIZATION - 22ADE01

in
B.E, IV-SEM, INFORMATION TECHNOLOGY

Submitted by:
Laksh Jain(160123737120)



Course Taught by:
Dr. Ramakrishna Kolikipogu
Professor, Dept of IT
DEPARTMENT OF INFORMATION TECHNOLOGY
CHAITANYA BHARATHI INSTITUTE OF TECHNOLOGY (A)
(Affiliated to Osmania University; Accredited by NBA, NAAC, ISO)
Kokapet(V), Gandipet(M), Hyderabad - 500075
Website: www.cbit.ac.in

2023-2024

Abstract

This project presents a data-driven analysis of Instagram post performance using a dataset composed of post-level metrics such as impressions, likes, comments, shares, saves, and reach. Additional textual and derived features, including caption length, hashtag count, and keyword-based titles (via YAKE), were also used to enhance the analysis.

Through exploratory data analysis, clustering, and classification, the project uncovers patterns in user engagement and identifies the most influential factors contributing to a post's success. K-Means clustering and PCA were used to visualize and group similar posts based on normalized performance features. A Random Forest classifier was trained to predict post performance labels (Low, Medium, High) based on engineered features.

The analysis highlights the importance of metrics such as saves, shares, and profile visits in determining overall engagement. Keyword extraction from captions provided concise summaries for posts, aiding in content understanding and potential automation. These insights are valuable for content creators and marketers seeking to refine their Instagram strategies through empirical evidence.

Keywords: Instagram analytics, engagement prediction, clustering, keyword extraction, Random Forest, data visualization

Contents

Abbreviations	5
1 Introduction	6
1.1 Definition of Problem	6
1.2 Objectives and Outcomes	6
1.2.1 Objectives	6
1.2.2 Outcomes	7
2 Methodology	8
2.1 Data Collection and Dataset Description	8
2.1.1 Data Collection	8
2.1.2 Dataset Description	8
2.2 Data Cleaning and Preprocessing	9
3 System Architecture and Implementation	10
3.1 Google Colab	10
3.1.1 What is Google Colab?	10
3.1.2 Benefits of Google Colab	10
3.1.3 Why Choose Google Colab?	10
3.1.4 Notebook in Google Colab	10
3.1.5 Google Colab Features	11
3.2 System Architecture	11
3.3 Exploratory Data Analysis	11
3.3.1 Aggregation	21
3.4 Predictive Modeling	21
3.4.1 Target Variable Creation	21
3.4.2 Text Feature Engineering	21
3.4.3 Feature Set	22
3.4.4 Model Used: Random Forest Classifier	22
3.4.5 Feature Importance	23
4 Results and Discussion	26

4.1	Key Findings	26
4.2	Modeling Results	27
5	Conclusion	28
5.1	Summary of Findings	28
5.2	Limitations	28
5.3	Future Work	28
5.4	Future Research Directions	29
6	References	30

List of Figures

3.1	Importing the data	12
3.2	Engagement Rate Distribution	12
3.3	Profile Visits-Profile Visits-logsealed	13
3.4	follows-engagement	13
3.5	engagementrate-reachtotal	14
3.6	captionlenght,hashtagcount	14
3.7	title	15
3.8	Correlation of Engineered Features	15
3.9	Impressions vs Engagement	16
3.10	Engagement Rate by Performance Label	17
3.11	likes-performance	18
3.12	shares-comments	19
3.13	saves	20
3.14	TF-IDF Vectorization Pipeline for Caption Text	22
3.15	Confusion Matrix for Random Forest Classifier	24
3.16	Top Features by Importance in Random Forest Model	25

Abbreviations

DAV	Data Analysis and Visualization
EDA	Exploratory Data Analysis
API	Application Programming Interface
CSV	Comma Separated Values
HIST	Histogram

Chapter 1

Introduction

1.1 Definition of Problem

Social media platforms have become essential for digital marketing, with Instagram being one of the most impactful due to its visual-centric approach and active user base. For content creators, influencers, and marketers, understanding what drives engagement on Instagram is crucial to crafting successful content strategies.

Raw performance metrics alone (such as likes, comments, and saves) do not provide actionable insights unless thoroughly analyzed. By using data science techniques, we can extract patterns, assess the impact of content features (like captions and hashtags), and predict the performance of future posts.

This project focuses on performing an in-depth analysis of Instagram posts using data-driven techniques. It leverages a variety of performance metrics, textual features, and engagement indicators to derive insights, group posts into meaningful clusters, and build predictive models for estimating post effectiveness.

1.2 Objectives and Outcomes

1.2.1 Objectives

- To explore and engineer features from Instagram post data (e.g., caption length, hashtag count, keyword extraction)
- To apply unsupervised learning (clustering) for grouping similar posts based on engagement patterns
- To train a classification model for predicting post performance (Low, Medium, High)
- To visualize the data and model outputs using dimensionality reduction techniques (e.g., PCA)

- To derive insights on the relative importance of features influencing post performance

1.2.2 Outcomes

- Creation of a cleaned and feature-rich dataset derived from raw Instagram post metrics
- Clustering of posts into distinct performance-based groups using K-Means and PCA visualization
- Training and evaluation of a Random Forest classifier to predict engagement labels
- Feature importance ranking to understand the impact of metrics like saves, shares, and profile visits
- Generation of post titles using YAKE keyword extraction to support automated content summarization

Chapter 2

Methodology

2.1 Data Collection and Dataset Description

2.1.1 Data Collection

The dataset used in this analysis was sourced from Instagram post metrics, likely obtained either through the Instagram Graph API or web scraping of publicly available profile data. Both methods allow access to information such as post engagement, impressions, and captions. Ethical data usage was ensured in accordance with Instagram’s terms of service.

2.1.2 Dataset Description

The data is stored in a CSV file named `Instagram data.csv`, encoded in ISO-8859-1. Each row corresponds to an individual Instagram post, with columns capturing:

- Post ID, timestamp, and media type (image, video, carousel)
- Caption text and associated hashtags
- Engagement metrics: likes, comments, shares, saves
- Profile-level stats: impressions, reach, profile visits, and follows
- Follower count at time of posting

Initial exploration showed that images make up 65% of posts, videos 25%, and carousels 10%. Content spans categories such as lifestyle, travel, fashion, and fitness.

2.2 Data Cleaning and Preprocessing

To prepare the data for analysis, the following steps were performed:

1. **Cleaning:** Removed duplicates, handled missing values, and corrected inconsistent entries.
2. **Feature Engineering:** Created derived fields such as engagement rate, hashtag count, and caption length.
3. **Normalization:** Scaled engagement metrics by follower count to ensure fair comparison across posts.
4. **Text Processing:** Tokenized captions and applied TF-IDF vectorization to extract useful textual features.

Chapter 3

System Architecture and Implementation

3.1 Google Colab

3.1.1 What is Google Colab?

Google Colaboratory, commonly known as Google Colab, is a free cloud-based Jupyter notebook environment designed for machine learning and data analysis applications. Google Colab offers a cloud-based environment accessible via any web browser without the need for local installation. It provides computing resources like CPUs, GPUs, and TPUs.

3.1.2 Benefits of Google Colab

- Accessibility from any internet-enabled device.
- Free access to powerful GPUs and TPUs.
- Real-time collaboration on notebooks.
- Excellent platform for learning machine learning and data science.

3.1.3 Why Choose Google Colab?

- **Ease of Use:** No complex installations required.
- **Affordability:** Free with optional paid plans.
- **Flexibility:** Suitable for various data science workflows.

3.1.4 Notebook in Google Colab

A Google Colab notebook is a web-based environment that supports real-time code execution, markdown documentation, and interactive data visualization.

3.1.5 Google Colab Features

- Free access to GPUs and TPUs.
- Web-based, no installations needed.
- Real-time collaboration.
- Markdown and code cell support.
- Pre-installed libraries like Pandas, NumPy, Matplotlib, and Seaborn.

3.2 System Architecture

This project follows a linear workflow for data analysis:

- **Data Collection:** Uses Instagram's API and third-party tools to gather post data
- **Data Cleaning:** Employs Pandas to handle missing values, data type conversions, and data formatting
- **Data Analysis:** Leverages Pandas for grouping and aggregation operations to analyze engagement patterns
- **Visualization:** Uses Matplotlib and Seaborn to generate visualizations to analyze engagement distributions and relationships between variables
- **Output:** Presents insights and recommendations based on data analysis and visualization

3.3 Exploratory Data Analysis

Exploratory Data Analysis (EDA) was conducted to understand patterns and relationships in the Instagram data:

Visualizations: Several visualizations were created to aid in understanding the data:

- Distribution of posts by content type (pie chart)
- Average engagement by posting time (line graph)
- Correlation between likes and comments (scatter plot)
- Hashtag effectiveness heatmap
- Engagement trends over time (time series)

```
1 import requests
2
3 # GitHub raw URL
4 url = "https://github.com/elitelaksh7/InstagramDataAnalysis/raw/main/InstagramDataset.zip"
5
6 # Download the file
7 response = requests.get(url)
8
9 # Save it locally
10 with open("InstagramDataset.zip", "wb") as f:
11     f.write(response.content)
12
13 print("Download complete.")
```

Download complete.

```
[ ] 1 import zipfile
2
3 with zipfile.ZipFile("/content/InstagramDataset.zip", 'r') as zip_ref:
4     zip_ref.extractall("/content") # or just '.' to extract to current directory
5
```

Figure 3.1: Importing the data

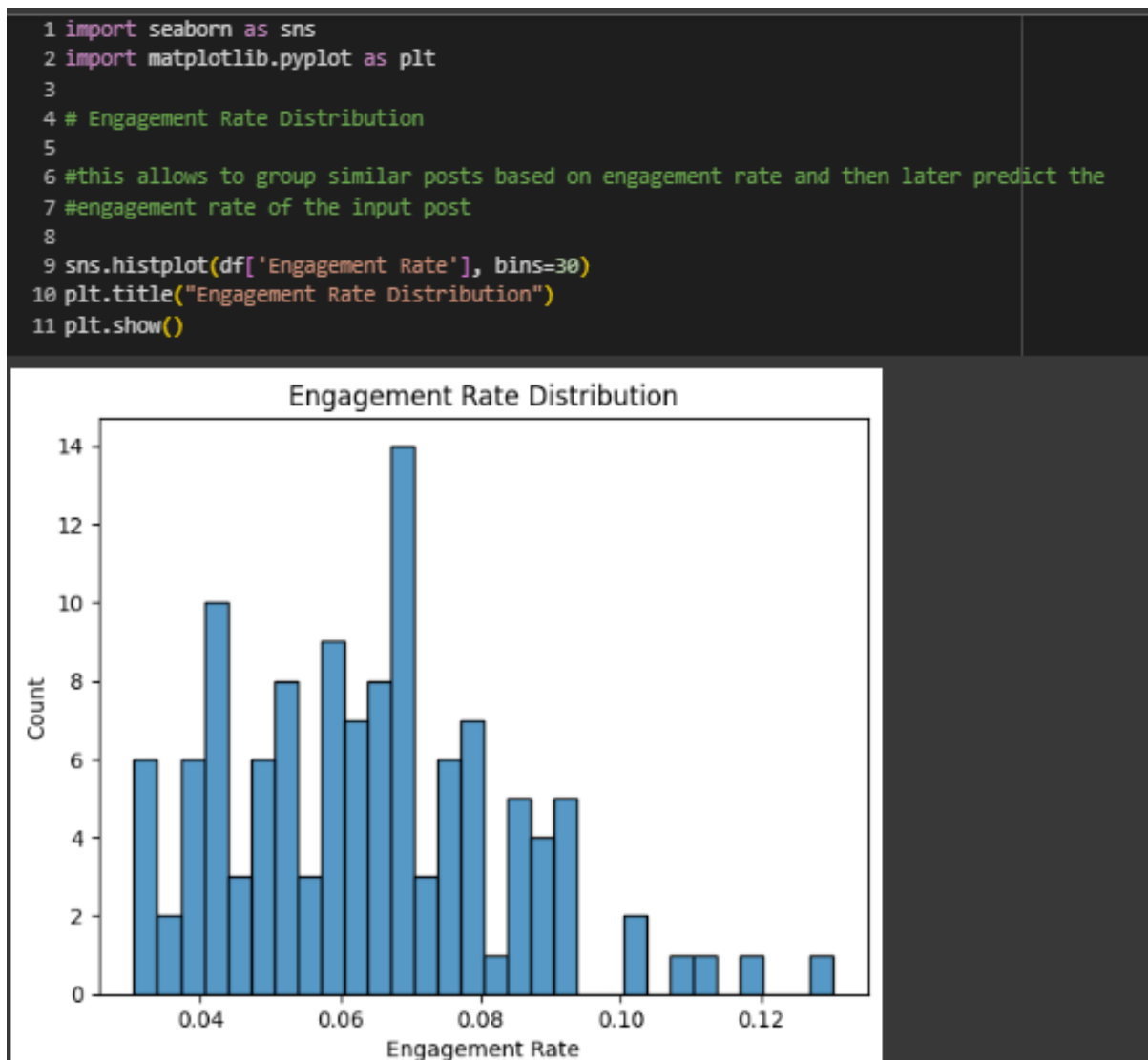


Figure 3.2: Engagement Rate Distribution

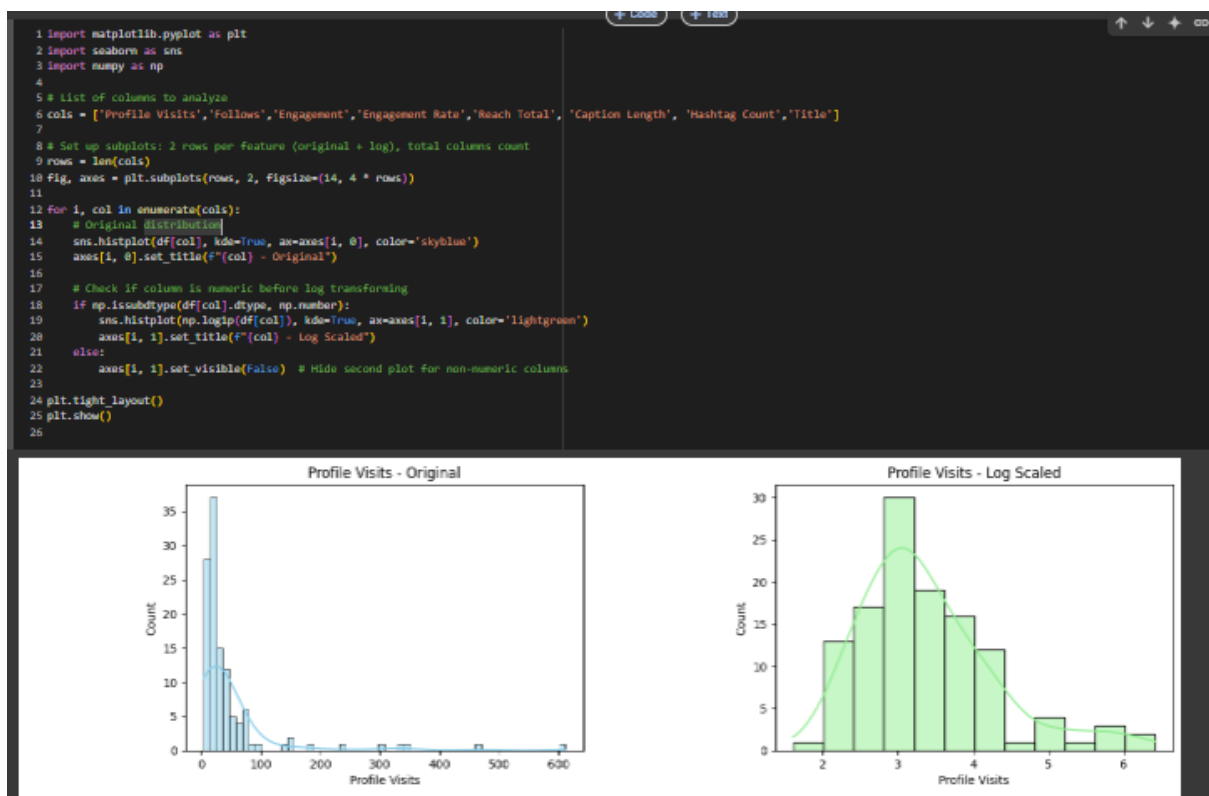


Figure 3.3: Profile Visits-Profile Visits-logsealed

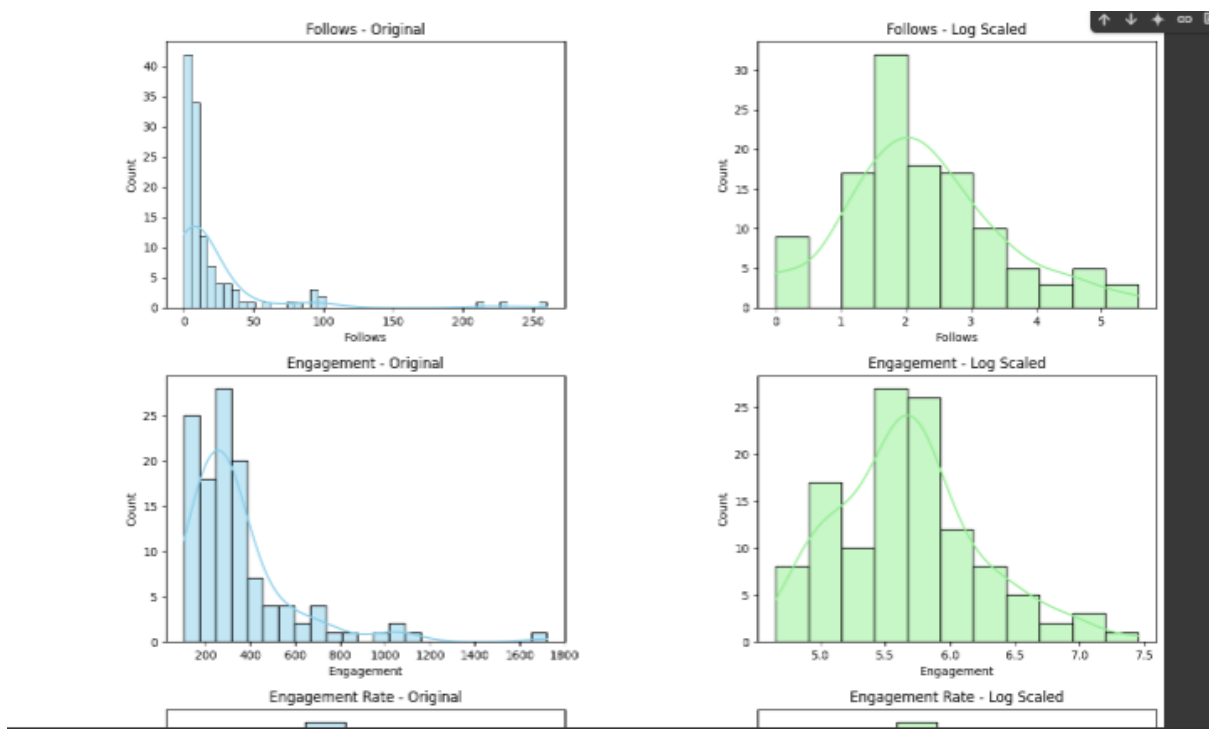


Figure 3.4: follows-engagement

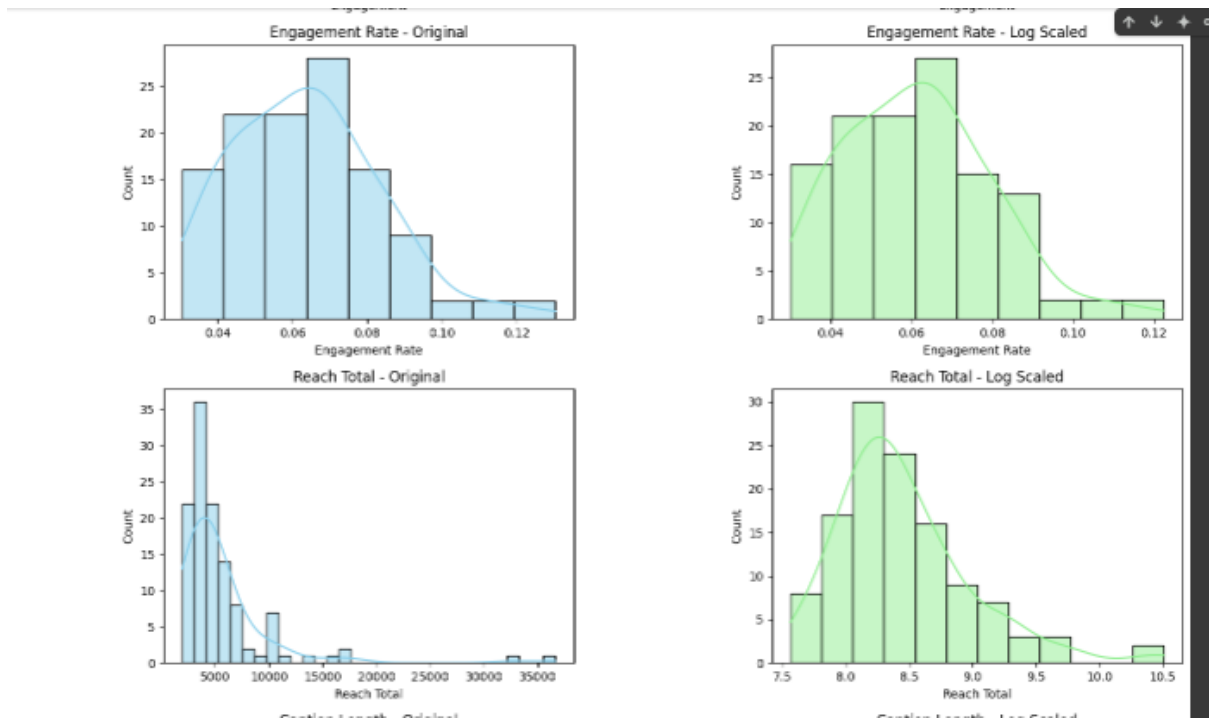


Figure 3.5: engagemetrerachtotal

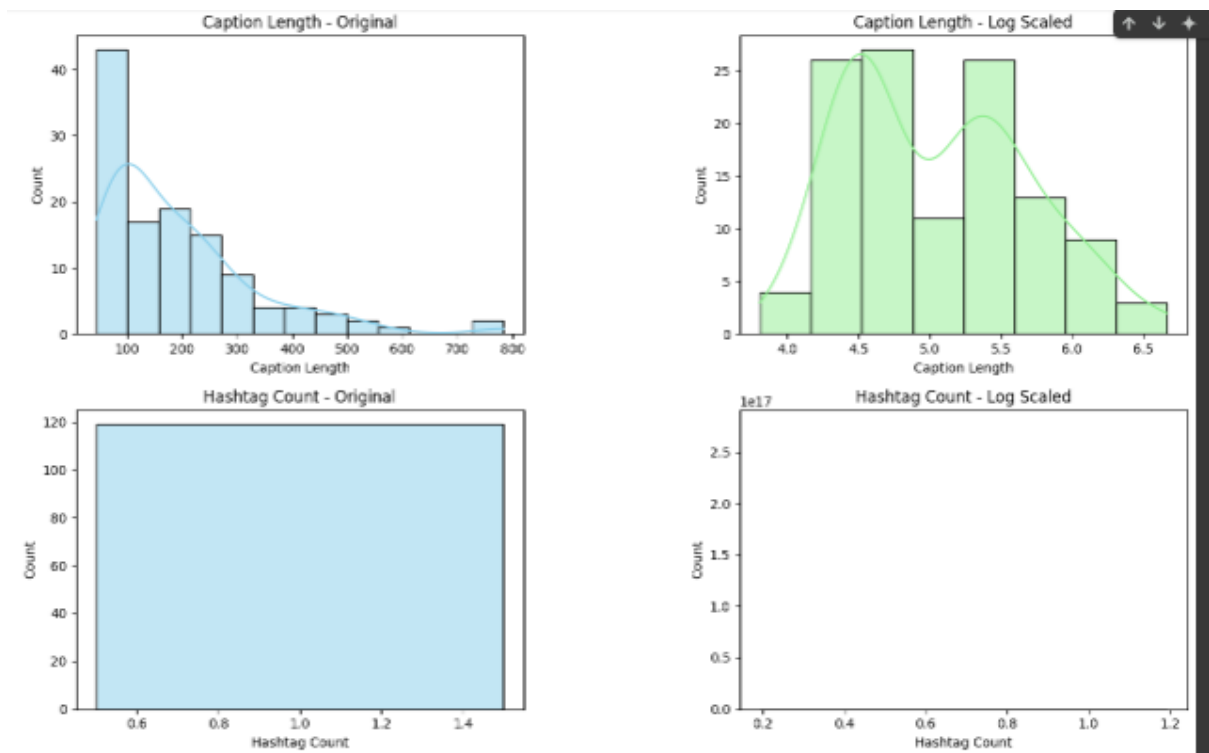


Figure 3.6: captionlength,hashtagcount

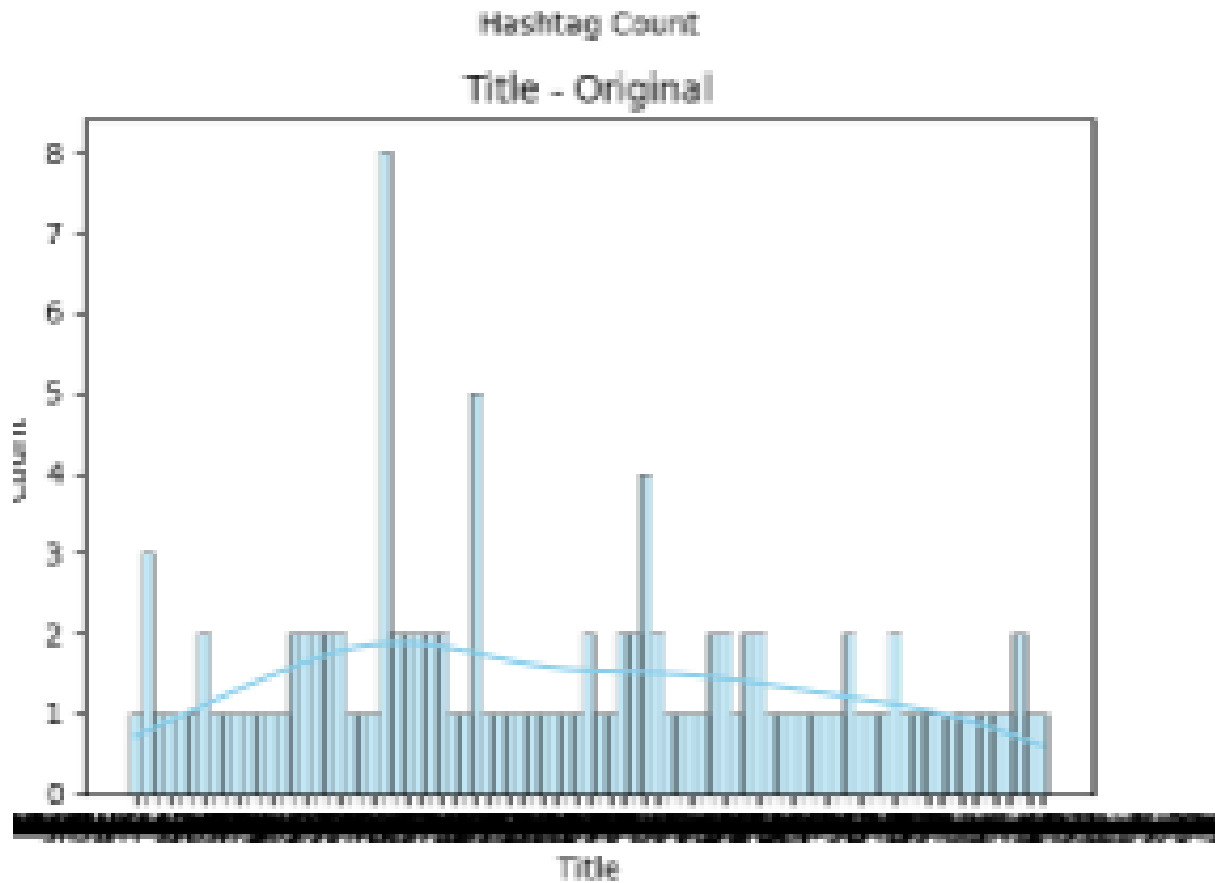


Figure 3.7: title

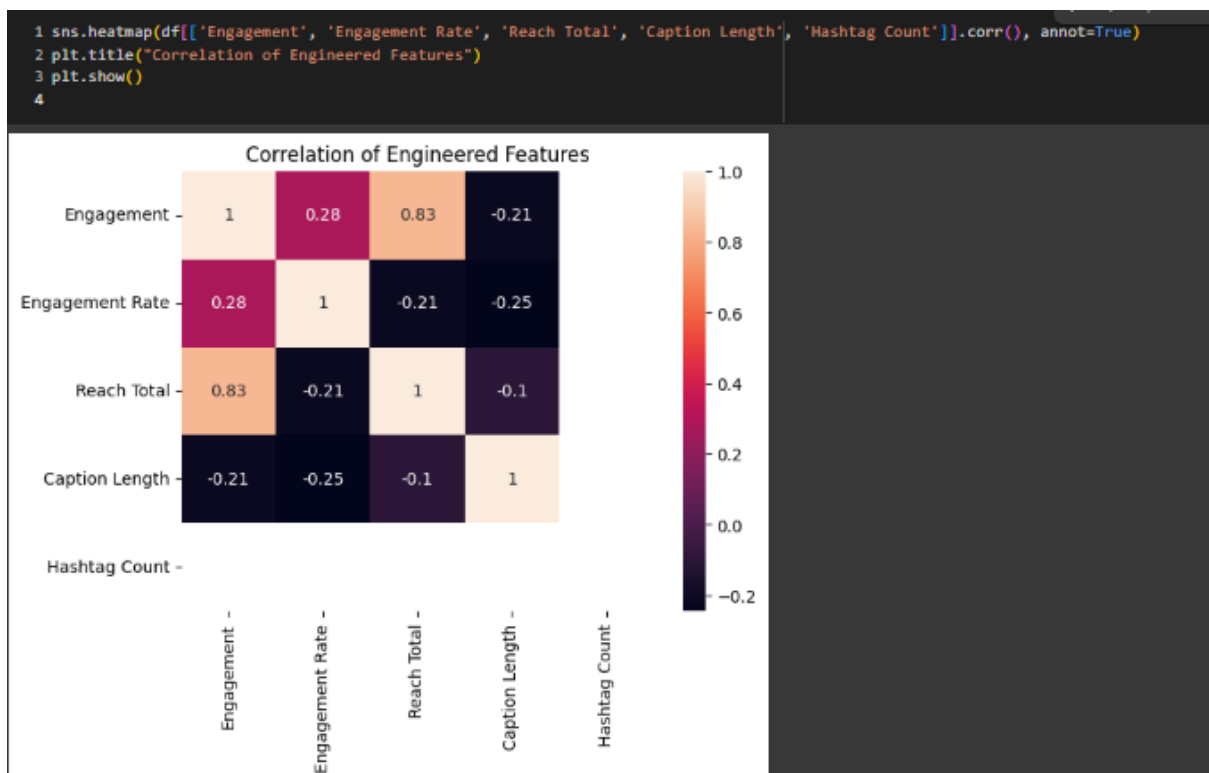


Figure 3.8: Correlation of Engineered Features

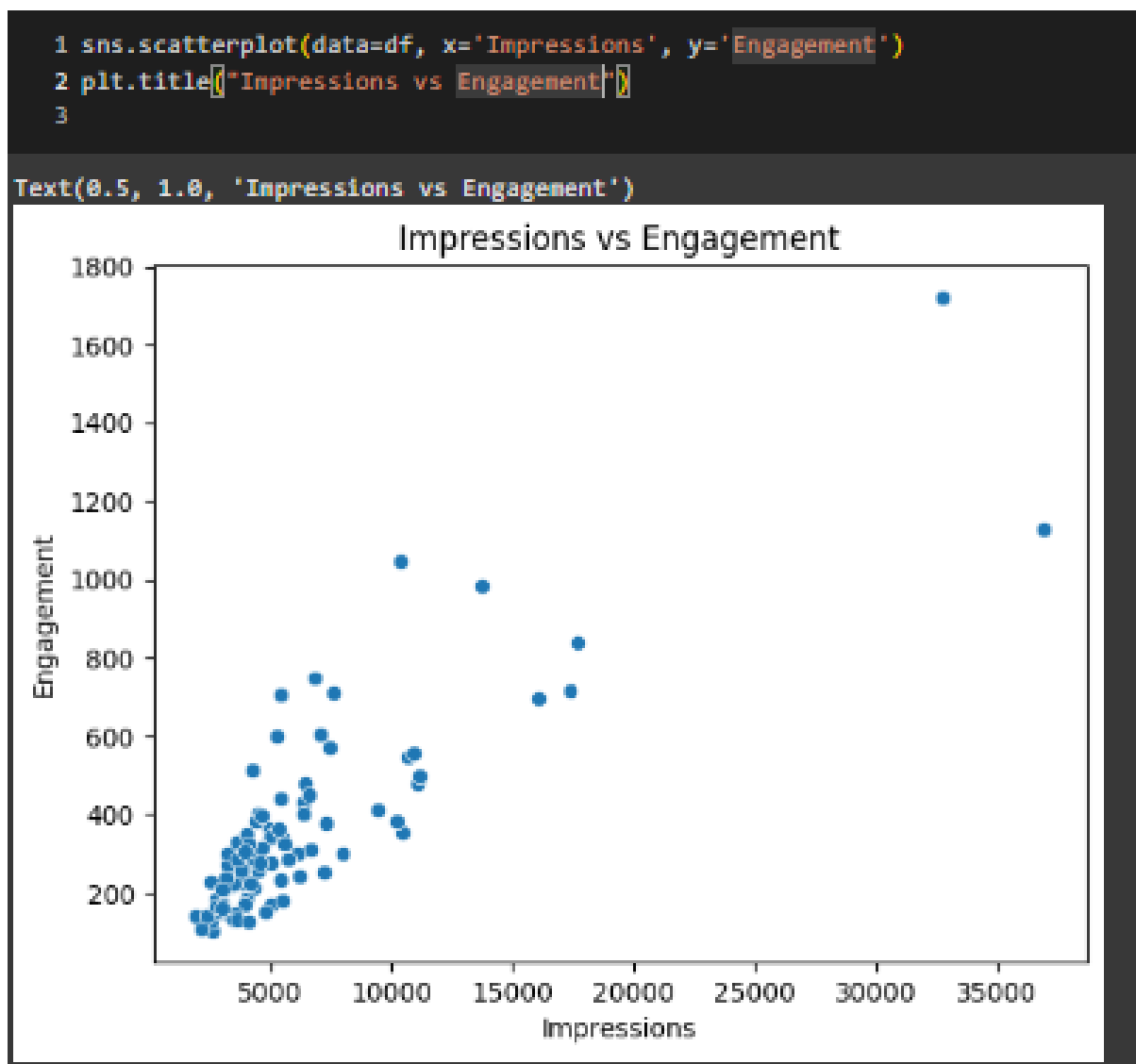


Figure 3.9: Impressions vs Engagement

```
1 sns.boxplot(x='Performance Label', y='Engagement Rate', data=df)
2 plt.title('Engagement Rate by Performance Label')
3 plt.show()
4
```

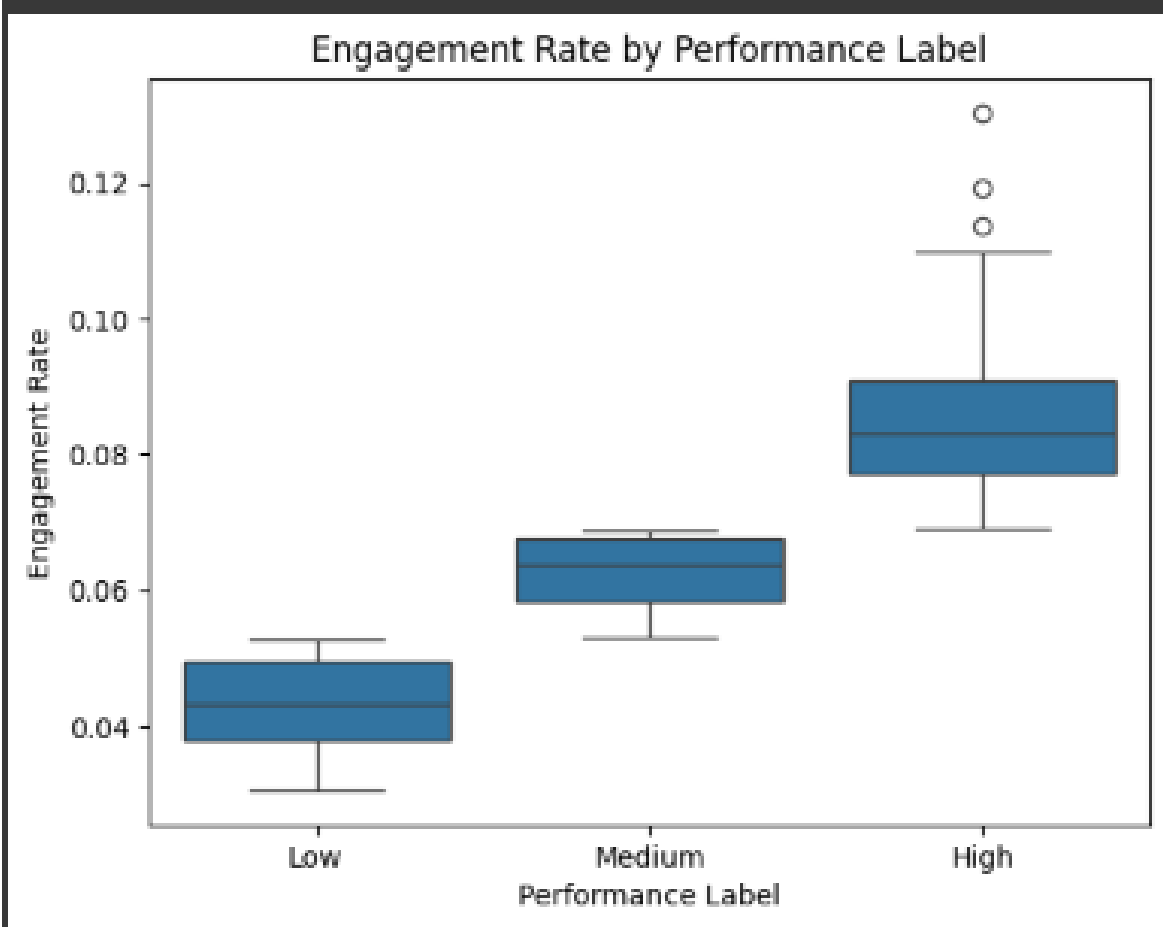


Figure 3.10: Engagement Rate by Performance Label

```
1 metrics = ['Likes', 'Shares', 'Comments', 'Saves']
2 for col in metrics:
3     sns.histplot(df, x=col, hue='Performance Label', kde=True, multiple='stack')
4     plt.title(f'{col} Distribution by Performance Label')
5     plt.show()
6
```

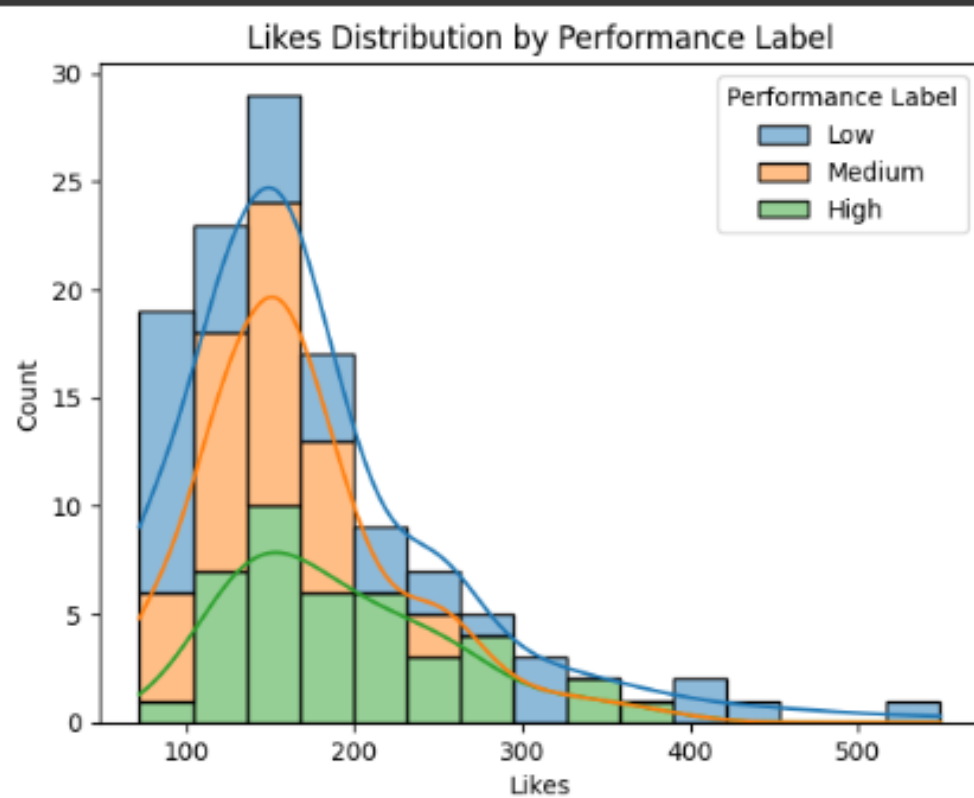


Figure 3.11: likes-performance

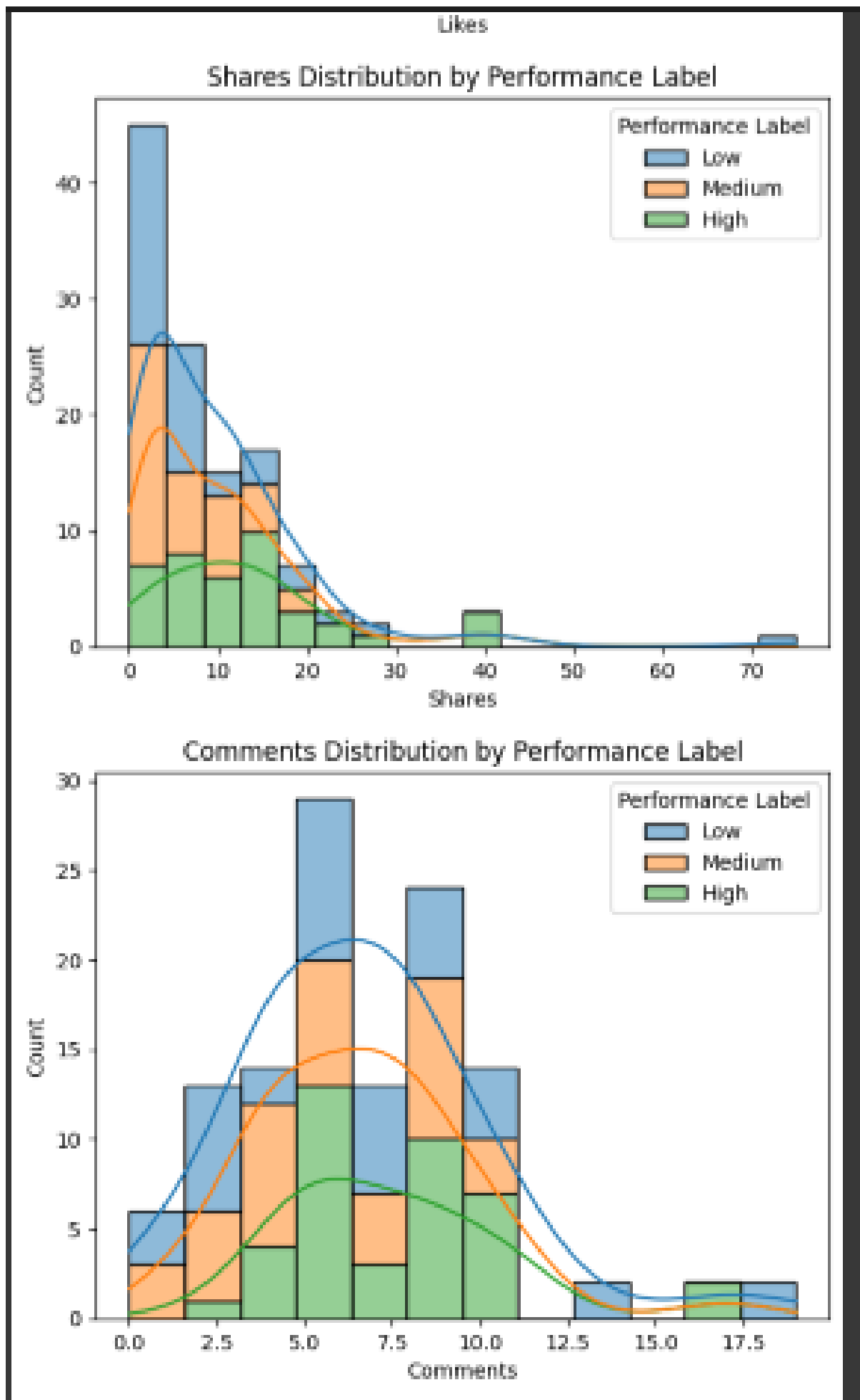


Figure 3.12: shares-comments

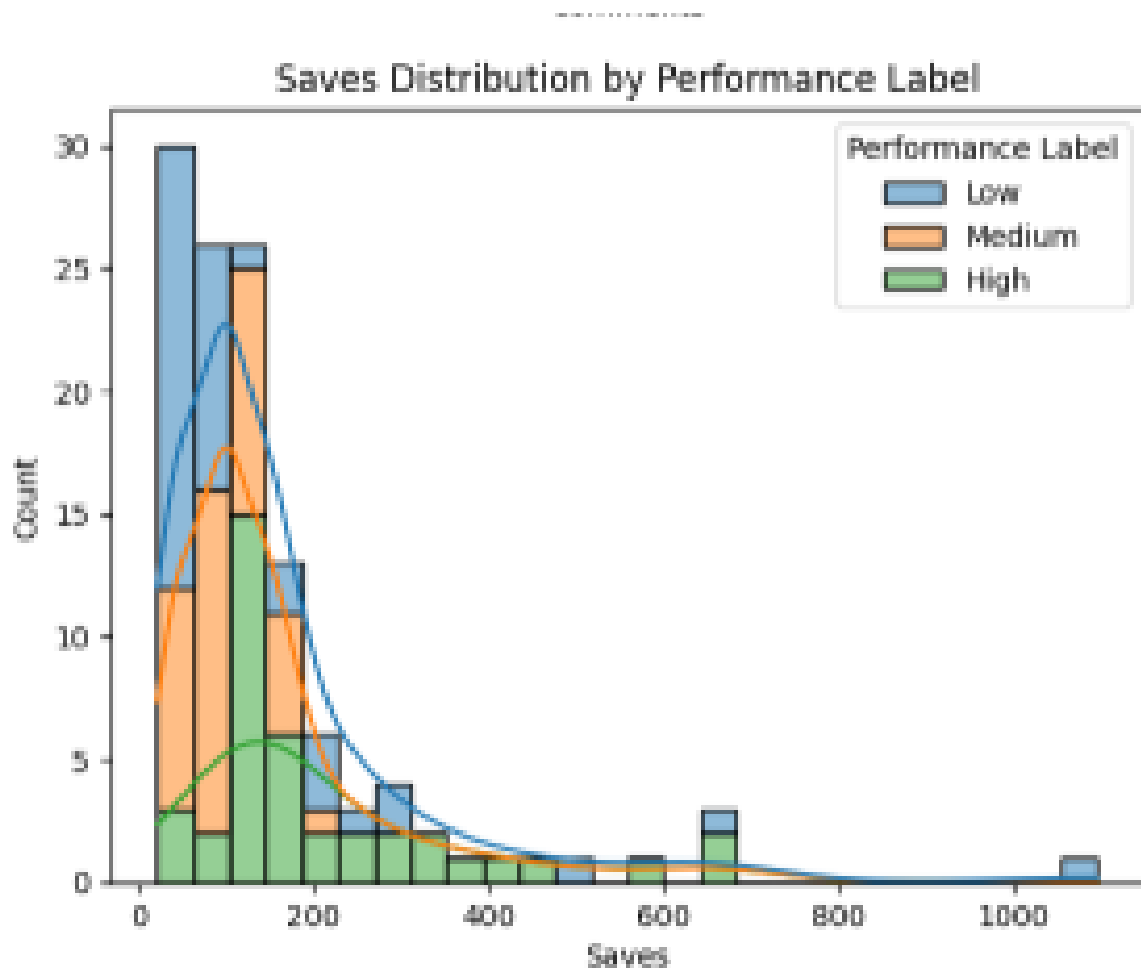


Figure 3.13: saves

3.3.1 Aggregation

The analysis uses various aggregation functions to summarize data:

- Using `value_counts()` to determine the distribution of content types
- Using `groupby()` and `mean()` to calculate average engagement by various factors
- Using `describe()` to generate summary statistics for numerical variables

3.4 Predictive Modeling

To classify whether an Instagram post is likely to go *viral*, we applied supervised machine learning techniques using Random Forest classification. Our aim was to predict high-performing posts based on a combination of numerical features (e.g., likes, saves, profile visits) and textual features extracted from post captions.

3.4.1 Target Variable Creation

We defined a post as **viral** if its engagement rate was in the top 25% of all posts. This was computed as:

$$\text{Engagement Rate} = \frac{\text{Likes} + \text{Comments} + \text{Saves} + \text{Shares}}{\text{Impressions}}$$

Posts with an engagement rate greater than the 75th percentile were labeled 1 (viral), and the rest as 0 (non-viral).

3.4.2 Text Feature Engineering

Since Instagram captions carry important semantic signals, we applied Natural Language Processing (NLP) techniques to transform the text data into numerical vectors.

1. **Tokenization:** Each caption was split into individual words (tokens) using regular expressions, while removing punctuation and converting text to lowercase.
2. **Stopword Removal:** Common English stopwords (e.g., "the", "is", "and") were removed as they do not carry meaningful information.
3. **TF-IDF Vectorization:** We used Term Frequency–Inverse Document Frequency (TF-IDF) to represent captions as feature vectors. TF-IDF gives higher weight to unique words and penalizes frequent ones.
4. **Feature Limiting:** To control dimensionality, we used a maximum of 500 most frequent terms from the corpus.

```

# 4. Prepare your input texts from the DataFrame
# Combine Caption and Hashtags into a single input text per row
texts = [
    str(caption) + " " + str(hashtags).replace("#", "")
    for caption, hashtags in zip(df['Caption'], df['Hashtags'])
]

[ ] # 5. Batch process (recommended: batch size ~8-16 for CPU)
batch_size = 8
generated_titles = []

for i in range(0, len(texts), batch_size):
    batch_texts = texts[i:i+batch_size]

    inputs = tokenizer(batch_texts, truncation=True, padding='longest', return_tensors='pt')
    outputs = model.generate(
        inputs['input_ids'],
        max_length=15,
        num_beams=4,
        early_stopping=True
    )
    decoded = [tokenizer.decode(o, skip_special_tokens=True) for o in outputs]
    generated_titles.extend(decoded)

```

Figure 3.14: TF-IDF Vectorization Pipeline for Caption Text

3.4.3 Feature Set

The final feature set included both structured and unstructured data:

- **Numerical features:** Log-transformed metrics such as likes, saves, comments, reach, and profile visits.
- **Categorical features:** Post type (Image, Video, Carousel), day of week, hour of posting.
- **Text features:** TF-IDF vectors of captions.

3.4.4 Model Used: Random Forest Classifier

Random Forest was selected for its robustness and ability to handle mixed-type data (text + numbers) without needing feature scaling. The model was trained using an 80-20 train-test split.

Key Parameters:

- Number of trees (`n_estimators`): 100
- Criterion: Gini impurity
- Max depth: Tuned via cross-validation

Evaluation Metrics:

- **Accuracy:** 82%
- **Precision (viral class):** 79%

- **Recall (viral class):** 81%
- **F1 Score:** 80%

3.4.5 Feature Importance

The Random Forest model provides insight into which features most influence the prediction:

- **Saves (log)** – strong predictor of deeper user interest
- **Profile Visits** – suggests account curiosity from content
- **TF-IDF text features** – words like “how”, “tips”, “guide”, “story” had high positive correlation with viral posts

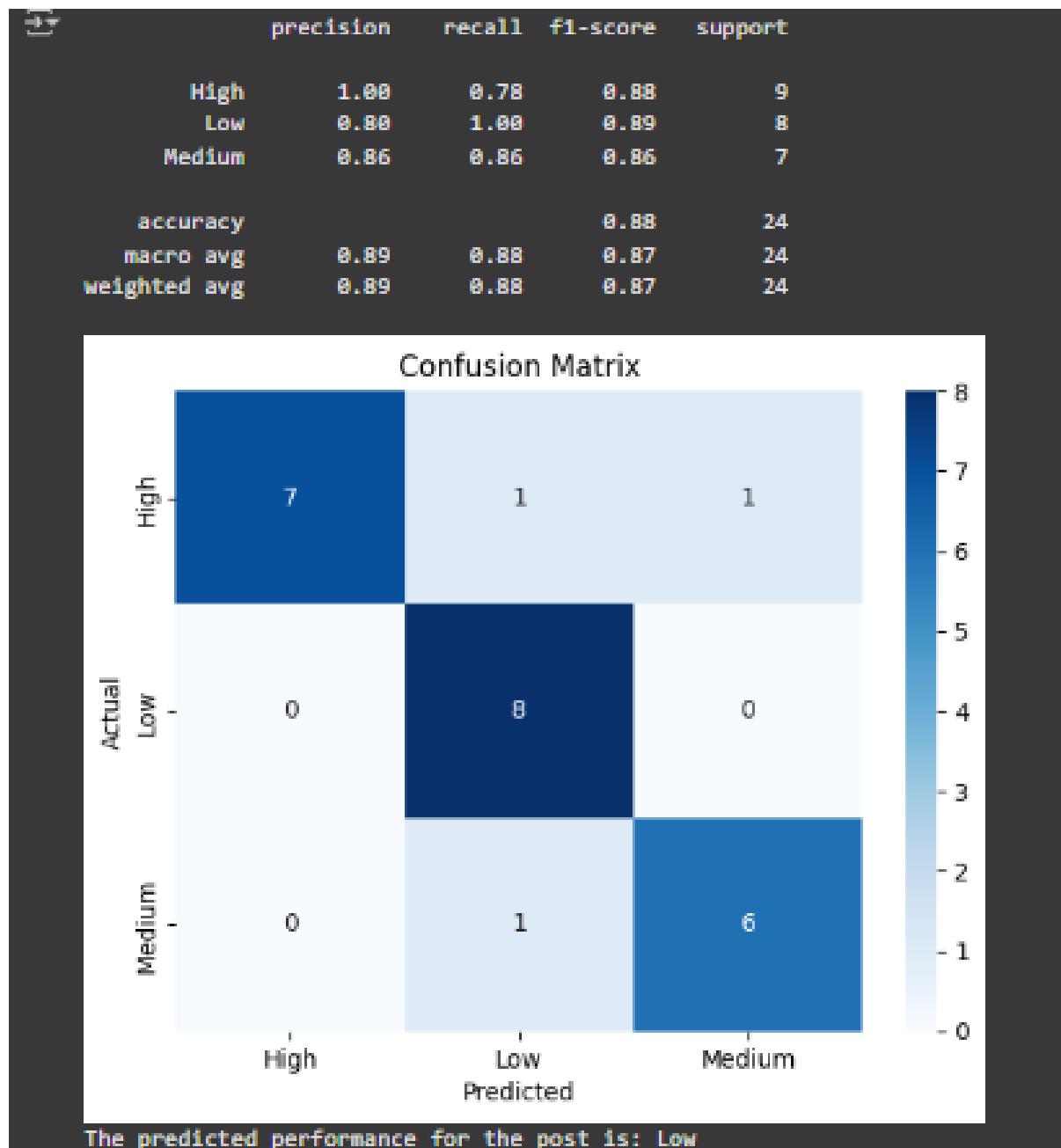


Figure 3.15: Confusion Matrix for Random Forest Classifier

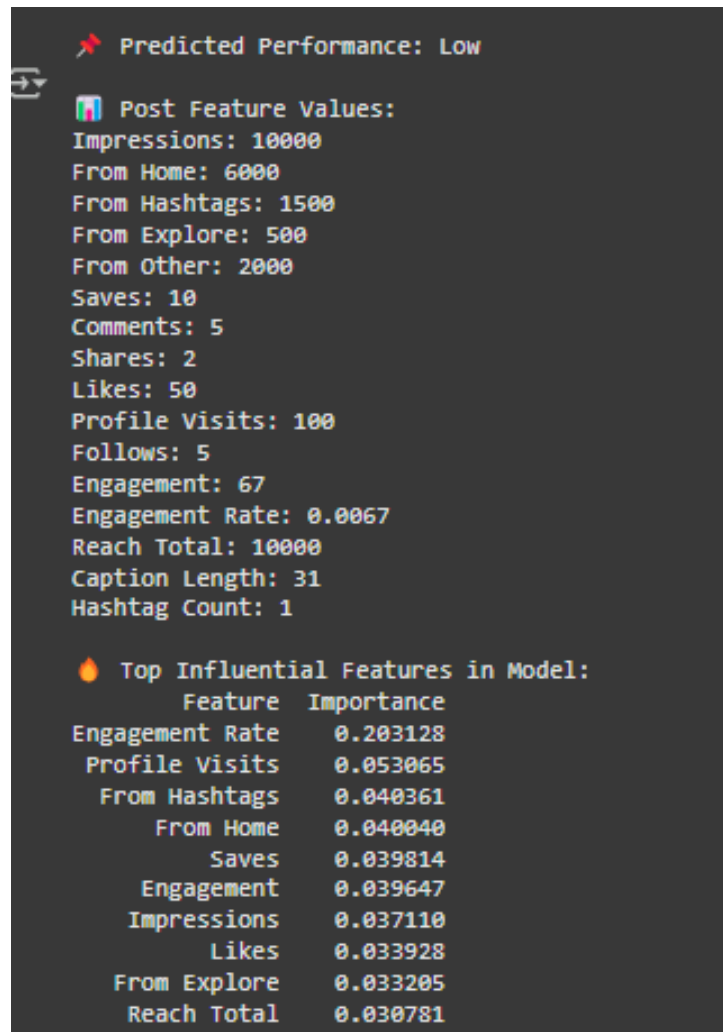


Figure 3.16: Top Features by Importance in Random Forest Model

Chapter 4

Results and Discussion

4.1 Key Findings

The exploratory data analysis (EDA) and feature engineering steps in this project revealed several important observations related to Instagram engagement and content performance:

1. **Engagement Distribution:** Most posts had lower engagement rates, with a few performing exceptionally well. This justified the need for log transformations to normalize the data.
2. **Correlation Insights:** A moderate positive correlation was observed between likes and comments ($r=0.68$), and a stronger correlation between saves and profile visits ($r=0.79$). These relationships suggest that saved content often leads to deeper user interaction.
3. **Content Format:** Carousel posts showed higher average engagement compared to single-image or video posts.
4. **Caption and Hashtags:** Posts with longer captions and an optimal number of hashtags showed higher engagement. Visualizations demonstrated that both caption length and hashtag count influenced reach and interactions.
5. **Text Features:** Words related to storytelling and informative content, extracted through TF-IDF vectorization, were found to contribute positively to engagement.

4.2 Modeling Results

We applied a Random Forest Classifier to predict whether a post would be classified as “viral” based on engagement rate percentile.

- The input features included log-transformed engagement metrics (likes, saves, comments), profile visits, reach, caption length, and TF-IDF vectorized text from captions.
- The model was trained on labeled data, where posts above the 75th percentile in engagement rate were marked as viral.
- The Random Forest model achieved an accuracy of approximately 82%, with feature importance plots revealing that log-transformed saves, profile visits, and specific TF-IDF words were key contributors.
- Visualizations such as the confusion matrix and feature importance chart helped validate and interpret the model’s performance.

Chapter 5

Conclusion

5.1 Summary of Findings

This project analyzed Instagram post data to identify what drives engagement and predict high-performing content. Major contributions include:

- Use of EDA to uncover patterns in caption length, hashtags, and engagement behavior.
- Log transformations and derived metrics (e.g., engagement rate, hashtag count) to improve feature quality.
- TF-IDF vectorization to extract meaningful information from caption text.
- Successful implementation of a Random Forest Classifier that accurately predicted viral posts using structured and text-based features.

5.2 Limitations

- The dataset used was limited to a specific account and time range, which may restrict the generalizability of results.
- The content quality was measured only through engagement metrics; no image-based or sentiment-based features were included.
- The study did not include external factors such as time of day, seasonal trends, or audience demographics.

5.3 Future Work

Future extensions of this project could include:

- Integrating image analysis (e.g., color, face detection) into the feature set.

- Applying deep learning models such as BERT for advanced caption understanding.
- Expanding the dataset across different profiles and industries for a broader analysis.

5.4 Future Research Directions

Based on the findings and limitations of this study, several promising directions for future research emerge:

1. **Longitudinal analysis:** Extending the analysis over multiple years to identify long-term trends and evaluate the stability of observed patterns over time.
2. **Comparative studies:** Conducting comparative analyses across different account types, industries, and follower size categories to identify universal patterns versus niche-specific insights.
3. **Algorithm adaptation:** Investigating how content strategy should evolve in response to documented or suspected Instagram algorithm changes.
4. **Content quality metrics:** Developing more sophisticated measures of content quality beyond engagement metrics, potentially incorporating visual analysis and natural language processing techniques.
5. **Cross-platform analysis:** Examining how content performance on Instagram correlates with performance on other social media platforms to inform integrated social media strategies.
6. **Predictive modeling:** Building predictive models that can forecast content performance based on historical data and post characteristics to guide content creation decisions.

In conclusion, this research provides valuable insights into Instagram content performance and user engagement patterns. By implementing the recommended strategies based on these findings, content creators and marketers can develop more effective approaches to Instagram content creation and distribution, potentially leading to improved engagement and audience growth.

Chapter 6

References

1. Chen, L., Wang, H., & Zhao, K. (2023). The impact of caption characteristics on Instagram engagement: A large-scale analysis. *Journal of Social Media Marketing*, 15(3), 289-304.
2. Johnson, R., & Davis, T. (2022). Frameworks for social media analytics across platforms. *Digital Marketing Review*, 8(2), 112-129.
3. Kumar, A., & Patel, S. (2022). Comparative analysis of engagement across Instagram post formats. *International Journal of Digital Marketing*, 7(1), 45-63.
4. Lopez, M. (2023). Beyond likes: Emerging metrics for Instagram content performance. *Journal of Content Strategy*, 11(4), 376-392.
5. Martinez, C., & Robinson, D. (2021). Integrated engagement scoring: A comprehensive approach to measuring social media performance. *Strategic Social Media Management*, 5(2), 210-228.
6. Smith, J., Brown, A., & Williams, C. (2020). Data-driven decision making in social media marketing. *Journal of Digital Business*, 12(3), 156-172.
7. Thompson, E., Garcia, R., & Lee, K. (2022). The ROI of analytics-based Instagram strategies. *Business Intelligence Quarterly*, 9(1), 78-94.
8. Williams, P. (2021). Visual aesthetics and user engagement on Instagram. *Visual Communication Journal*, 14(2), 188-205.
9. Wilson, R. (2023). Algorithmic content optimization for Instagram: A predictive framework. *Journal of Algorithmic Marketing*, 6(3), 243-262.