# Recipe Reviews Sentiment Analysis

# Chapter 1

# Introduction

## 1.1 Business Context

This project focuses on analyzing user feedback and reviews on recipes using Apache Spark to uncover patterns in user engagement and sentiment on a recipe-sharing platform. The purpose is to enhance recipe recommendations and improve customer satisfaction, leveraging the large-scale data processing capabilities of Spark.

## 1.2 Dataset Overview

The dataset is sourced from the UCI Machine Learning Repository. It consists of user-generated feedback from an online platform called Food.com that ranks recipes. The data, likely collected through web scraping, includes extensive user interactions such as comments, ratings, and reputations, making it an ideal source for detailed behavioral analysis.

- **File Size**: 5.8 MB

- **Format**: CSV (Comma-Separated Values)

- **Number of Instances**: 18,182 rows

- **Number of Features**: 15 columns

### 1.2.1 Key Attributes of the Dataset

- **Recipe Details**: `recipe_name`, `recipe_number` (rank on the top 100 list), `recipe_code` (unique ID).

- **User Details**: `user_id`, `user_name`, `user_reputation`.

- **Review Metadata**: `comment_id`, `created_at` (Unix timestamp), `reply_count`, `thumbs_up`, `thumbs_down`, `stars` (rating from 1 to 5, with 0 indicating no rating).

- **Textual Data**: The `text` column contains the full content of the user comment, which may include reviews or feedback on the recipe.

## 1.3 Purpose

The project focuses on exploring a rich dataset of recipe reviews and user interactions sourced from an online platform, making it highly relevant for a range of practical applications. The primary purpose is to leverage this data for sentiment analysis, user behavior modeling, and trend identification using advanced data processing techniques. By analyzing attributes like `thumbs_up`, `rating`, `user_reputation`, and textual feedback, the project aims to uncover user preferences and engagement patterns. Additionally, the dataset supports the development of robust recommendation systems tailored to user tastes and preferences. Researchers can utilize this dataset for machine learning tasks, including text classification and user interaction modeling. The project demonstrates the use of PySpark and NLP techniques to extract insights, highlighting its value in improving user experience, understanding content engagement, and developing data-driven strategies in the domain of online recipe platforms.

# Chapter 2

# Data Preparation and Preprocessing

## 2.1 Initialize Spark Session

This step initializes a Spark session, which serves as the foundation for running a PySpark application.

```python
import findspark
findspark.init()

from pyspark.sql import SparkSession

# Initialize Spark Session
spark = SparkSession.builder \
    .master("local[4]") \
    .appName("Recipe Reviews EDA") \
    .enableHiveSupport() \
    .getOrCreate()

# Load the data
sc = spark.sparkContext

# Extract and print the port number of Spark's Web UI
spark_session_port = spark.sparkContext.uiWebUrl.split(":")[-1]
print("Spark Session WebUI Port: " + spark_session_port)

# Set logging level to ERROR
sc.setLogLevel("ERROR")

# Show available databases
databases_df = spark.sql("show databases")
databases_df.show()

# Show available tables
tables_df = spark.sql("show tables")
tables_df.show()
```

Listing 2.1: Initializing Spark Session

## 2.2 Load and Explore Dataset

The dataset, stored as a CSV file, is loaded into a Spark DataFrame. The `header=True` option ensures the first row is treated as column names, and `inferSchema=True` allows

Spark to automatically determine the data types of each column.

```
1 df = spark.read.csv('Recipe Reviews and User Feedback Dataset.csv',
    header=True, inferSchema=True)
2
3 # Display the first 5 rows of the dataframe
4 df.show(5)
5
6 # Display the schema of the DataFrame
7 df.printSchema()
```

Listing 2.2: Loading and Exploring Dataset

## 2.3   Data Type Casting

To prepare the data for analysis, specific columns are cast into appropriate data types. For example, `recipe_number` and `thumbs_up` are cast to `IntegerType` to ensure they can participate in arithmetic operations, and `created_at` is cast to `TimestampType` for time-based computations. Casting ensures data consistency and avoids errors in subsequent steps.

```
1 from pyspark.sql.types import IntegerType, TimestampType
2 from pyspark.sql.functions import col
3
4 # Cast columns to the appropriate types
5 df = df.withColumn("recipe_number", col("recipe_number").cast(
    IntegerType())) \
6        .withColumn("recipe_code", col("recipe_code").cast(IntegerType()
    )) \
7        .withColumn("user_reputation", col("user_reputation").cast(
    IntegerType())) \
8        .withColumn("created_at", col("created_at").cast(TimestampType()
    )) \
9        .withColumn("reply_count", col("reply_count").cast(IntegerType()
    )) \
10        .withColumn("thumbs_up", col("thumbs_up").cast(IntegerType())) \
11        .withColumn("thumbs_down", col("thumbs_down").cast(IntegerType()
    )) \
12        .withColumn("stars", col("stars").cast(IntegerType())) \
13        .withColumn("best_score", col("best_score").cast(IntegerType()))
14
15 # Verify the schema
16 df.printSchema()
17
18 # Rename columns for better clarity if needed
19 df = df.withColumnRenamed("_c0", "id") \
20        .withColumnRenamed("created_at", "review_timestamp") \
21        .withColumnRenamed("stars", "rating") \
22        .withColumnRenamed("text", "review_comments")
23
24 from pyspark.sql.functions import year, month, dayofmonth, hour,
    dayofweek
25
26 # Extract year, month, day, hour, and weekday
27 df = df.withColumn("year", year(df["review_timestamp"])) \
28        .withColumn("month", month(df["review_timestamp"])) \
29        .withColumn("day", dayofmonth(df["review_timestamp"])) \
```

```
30        .withColumn("hour", hour(df["review_timestamp"])) \
31        .withColumn("day_of_week", dayofweek(df["review_timestamp"]))  #
      Monday = 1, Sunday = 7
```

Listing 2.3: Data Type Casting

## 2.4   Data Cleaning

This step checks for missing data in all columns and fills non-critical fields, such as
`rating` and `user_reputation`, with default values like 0. To maintain data quality, rows
with missing values in important columns like `recipe_number` and `review_comments` are
removed, ensuring the dataset is clean and ready for analysis.

```
1 from pyspark.sql.functions import when, count, isnan
2
3 # Check for missing values in each column
4 df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns
    ]).show()
5
6 # Handle missing values
7 df = df.na.fill({ "rating" : 0, "user_reputation": 0, "thumbs_up":0, "
    thumbs_down":0})
8
9 # Drop rows with nulls in critical columns
10 df = df.na.drop(subset=["recipe_number", "recipe_code", "user_id", "
    comment_id", "review_comments"])
11
12 # Create a temporary view for SQL operations
13 df.createOrReplaceTempView("recipe_reviews")
14
15 # Descriptive Statistics
16 desc_stats_df = df.describe(["rating", "user_reputation", "thumbs_up",
    "thumbs_down", "reply_count"])
17 desc_stats_df.show()
```

Listing 2.4: Data Cleaning

## 2.5   Correlation Analysis

This step analyzes the relationship between `user_reputation`, `thumbs_up`, and `rating`
to identify potential dependencies. A heatmap is used to visualize the correlation values,
making it easier to understand the strength and direction of relationships.

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Correlation analysis between user_reputation, thumbs_up, and rating
5 correlation_df = df.select("user_reputation", "thumbs_up", "rating").
    toPandas()
6
7 # Compute correlations between user_reputation, thumbs_up, and rating
8 correlations = correlation_df.corr()
9 print(correlations)
10
11 # Plot the correlation heatmap
```
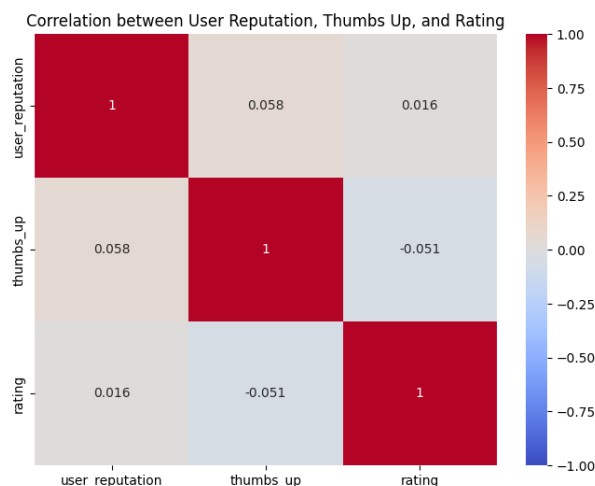
Figure 2.1: Correlation between User Reputation, Thumbs Up, and Rating

```
12  plt.figure(figsize=(8, 6))
13  sns.heatmap(correlations, annot=True, cmap="coolwarm", vmin=-1, vmax=1)
14  plt.title("Correlation between User Reputation, Thumbs Up, and Rating")
15  plt.show()
```

Listing 2.5: Correlation Analysis

## 2.6 Visualization of Top Recipes

The following queries and visualizations identify the top recipes based on average ratings and thumbs-up votes.

### 2.6.1 Top 20 Recipes by Average Rating

```
1   # Query to see distribution of star ratings
2   top_recipes_df = spark.sql("""
3       SELECT recipe_name, AVG(rating) AS average_rating
4       FROM recipe_reviews
5       GROUP BY recipe_name
6       ORDER BY average_rating DESC
7       LIMIT 20
8   """)
9
10  top_recipes_df.show()
11
12  # Convert to Pandas for visualization
13  top_recipes_pd = top_recipes_df.toPandas()
14
15  plt.figure(figsize=(10, 6))
16  plt.barh(top_recipes_pd['recipe_name'], top_recipes_pd['average_rating'
        ], color='orange')
17  plt.xlabel('Average Rating')
18  plt.title('Top 20 Recipes with the Highest Average Ratings')
19  plt.gca().invert_yaxis()  # To have the highest rating at the top
20  plt.show()
```

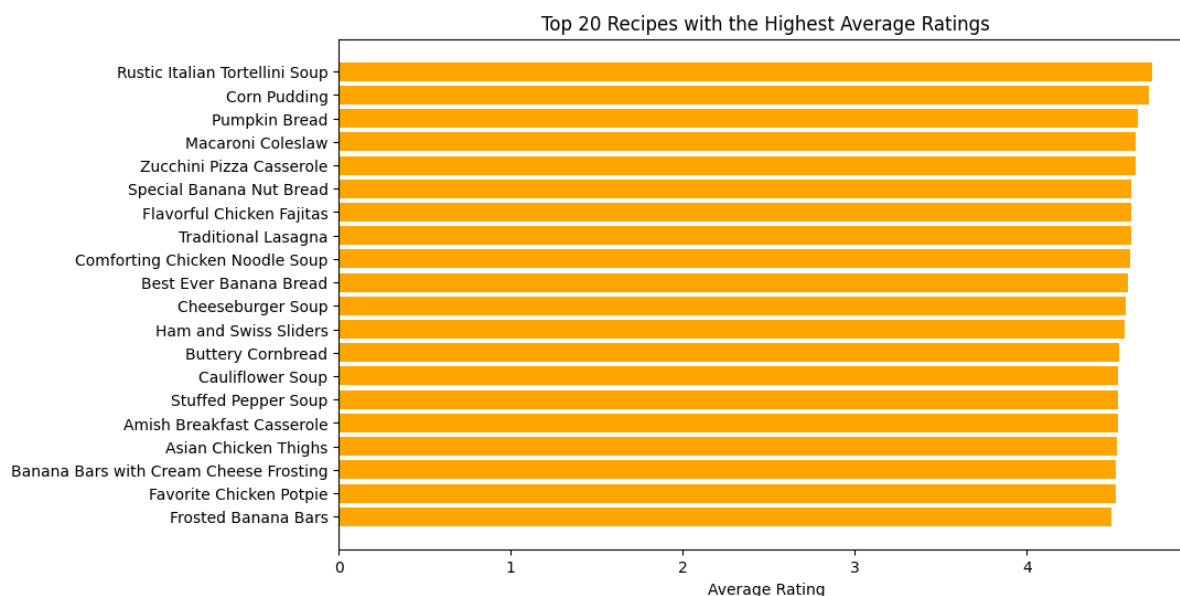Listing 2.6: Top 20 Recipes by Average Rating

Figure 2.2: Top 20 Recipes with the Highest Average Ratings

## 2.6.2   Top 20 Recipes by Thumbs Up

```
thumbs_df = spark.sql("""
    SELECT
        recipe_name,
        SUM(thumbs_up) as total_thumbs_up,
        SUM(thumbs_down) as total_thumbs_down,
        (SUM(thumbs_up) + SUM(thumbs_down)) as total_thumbs
    FROM recipe_reviews
    GROUP BY recipe_name
    ORDER BY total_thumbs_up DESC
    LIMIT 20
""")

thumbs_df.show()

# Visualization
engaged_recipes_df = spark.sql("""
SELECT recipe_name,
    SUM(thumbs_up) AS total_thumbs_up,
    SUM(thumbs_down) AS total_thumbs_down,
    AVG(rating) AS avg_rating
    FROM recipe_reviews
    GROUP BY recipe_name
    ORDER BY total_thumbs_up DESC
    LIMIT 20
""")

# Convert to Pandas for visualization
engaged_recipes_pd = engaged_recipes_df.toPandas()

# Plot
plt.figure(figsize=(12, 6))
plt.bar(engaged_recipes_pd['recipe_name'], engaged_recipes_pd[
    'total_thumbs_up'], label='Thumbs Up', color='green')
```

Figure 2.3: Top 20 Most Engaged Recipe

```
33 plt.bar(engaged_recipes_pd['recipe_name'], engaged_recipes_pd['
      total_thumbs_down'], label='Thumbs Down', color='red', bottom=
      engaged_recipes_pd['total_thumbs_up'])
34 plt.xlabel('Recipe Name')
35 plt.ylabel('Engagement Count')
36 plt.title('Top 20 Most Engaged Recipes')
37 plt.xticks(rotation=90)
38 plt.legend()
39 plt.show()
```

Listing 2.7: Top 20 Recipes by Thumbs Up

## 2.6.3   Review Length vs Rating

```
1 from pyspark.sql.functions import length
2
3 # Length of review text vs rating
4 df = df.withColumn("review_length", length(col("review_comments")))
5 sns.boxenplot(data=df.toPandas(), x='rating', y='review_length', color=
      'brown')
6 plt.title('Review Length vs Rating')
7 plt.xlabel('Review Length (Characters)')
```

Figure 2.4: Review Length vs Rating

```
8 plt.ylabel('Rating')
9 plt.show()
```

Listing 2.8: Review Length vs Rating

## 2.6.4 Ratings Over Time

```
1 # Ratings over time
2 trends_df = spark.sql("""
3     SELECT year, month, AVG(rating) AS avg_rating
4     FROM recipe_reviews
5     GROUP BY year, month
6     ORDER BY year, month
7 """)
8
9 # Convert to Pandas for visualization
10 temporal_trends_pd = trends_df.toPandas()
11
12 # Plot temporal trends for average ratings over months
13 plt.figure(figsize=(12, 6))
14 sns.lineplot(x="month", y="avg_rating", data=temporal_trends_pd, marker
    ="o", hue="year")
15 plt.title("Average Rating Over Time by Year and Month")
16 plt.xlabel("Month")
17 plt.ylabel("Average Rating")
18 plt.show()
```

Listing 2.9: Ratings Over Time

## 2.6.5 Sentiment Analysis Using TextBlob

```
1 from textblob import TextBlob
2 import matplotlib.pyplot as plt
3
4 # Convert the DataFrame to an RDD
5 comments_rdd = df.select("review_comments", "rating").dropna().rdd
6
```

Figure 2.5: Average Rating Over Time by Year and Month

```
7  # Filter out ratings with a value of 0
8  comments_rdd_filtered = comments_rdd.filter(lambda row: row["rating"]
       != 0)
9
10 # Map each comment to its sentiment score
11 sentiment_rdd = comments_rdd.map(lambda row: (row["rating"], TextBlob(
       row["review_comments"]).sentiment.polarity))
12
13 # Map to (rating, (sentiment_score, 1)) for counting and summing scores
14 sentiment_count_rdd = sentiment_rdd.mapValues(lambda sentiment: (
       sentiment, 1))
15
16 # Reduce by key to sum up the sentiment scores and counts for each
       rating
17 sum_counts_by_rating = sentiment_count_rdd.reduceByKey(lambda x, y: (x
       [0] + y[0], x[1] + y[1]))
18
19 # Calculate the average sentiment for each rating
20 avg_sentiment_by_rating = sum_counts_by_rating.mapValues(lambda x: x[0]
       / x[1]).collect()
21
22 # Convert to a sorted list for plotting
23 avg_sentiment_by_rating = sorted(avg_sentiment_by_rating, key=lambda x:
       x[0])
24
25 # Separate ratings and sentiment values for plotting
26 ratings, avg_sentiments = zip(*avg_sentiment_by_rating)
27
28 # Plot the average sentiment by rating
29 plt.figure(figsize=(10, 6))
30 plt.bar(ratings, avg_sentiments, color='coral')
31 plt.xlabel("Rating")
32 plt.ylabel("Average Sentiment Score")
33 plt.title("Average Sentiment Score by Rating")
34 plt.show()
```
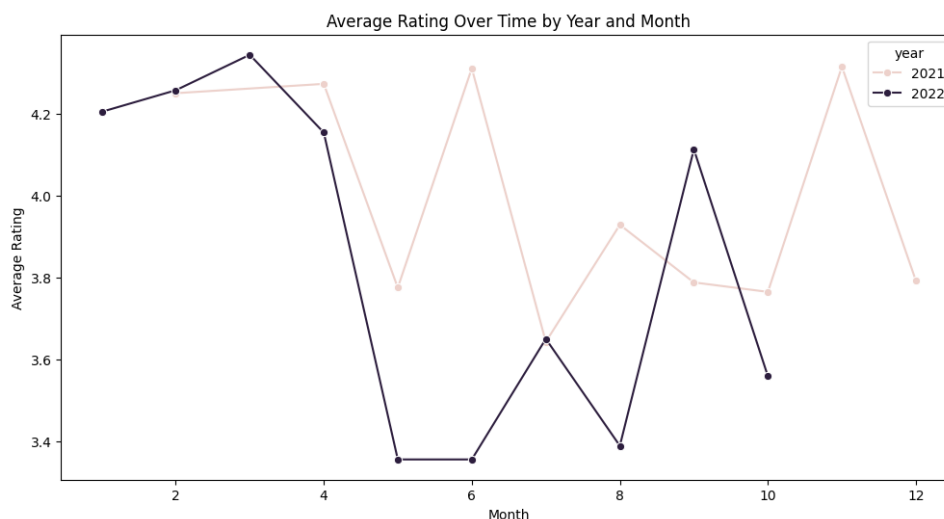
Listing 2.10: Sentiment Analysis Using TextBlob

Figure 2.6: Average Sentiment Score by Rating

## 2.7   Noun and Adjective Analysis

```python
from pyspark.sql import SparkSession
from tqdm import tqdm
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer, PorterStemmer
import re
import nltk

# Download required NLTK resources
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')

# Initialize NLTK tools
stop_words = set(stopwords.words("english"))

# Add custom stopwords for unwanted words
custom_stopwords = {"good", "great", "easy", "much", "way", "next", "
    little", "time", "recipe", "delicious"}
stop_words.update(custom_stopwords)

lemmatizer = WordNetLemmatizer()

# Function to clean and preprocess text
def preprocess_text(text):
    # Convert to lowercase
    text = text.lower()

    # Remove special characters, numbers, and single characters
    text = re.sub(r'\b\w{1}\b', '', text)   # Remove single letters
    text = re.sub(r'\d+', '', text)          # Remove numbers
    text = re.sub(r'[^\w\s]', '', text)       # Remove special
    characters
```

```python
31
32     # Tokenize text and filter by POS tags for nouns only
33     blob = TextBlob(text)
34     tokens = [word for word, pos in blob.tags if pos in ('NN', 'NNS')]
35
36     # Remove stopwords and lemmatize each word
37     tokens = [lemmatizer.lemmatize(word) for word in tokens if word not
       in stop_words]
38
39     return tokens
40
41 # Initialize Spark session
42 spark = SparkSession.builder \
43     .appName("Noun and Adjective Count by Rating") \
44     .getOrCreate()
45
46 # Convert DataFrame to an RDD with (rating, comment) pairs
47 comments_rdd = df.select("review_comments", "rating").dropna().rdd
48
49 # Preprocess and extract features for each comment
50 processed_comments_rdd = comments_rdd.map(lambda row: (row["rating"],
       preprocess_text(row["review_comments"])))
51
52 # Flatten the list of features and map each to (rating, feature) pairs
53 flattened_features_rdd = processed_comments_rdd.flatMap(lambda x: [(x
       [0], feature) for feature in x[1]])
54
55 # Count each noun and adjective for each rating
56 feature_counts_rdd = flattened_features_rdd.map(lambda x: ((x[0], x[1])
       , 1)).reduceByKey(lambda a, b: a + b)
57
58 # Transform back to (rating, (feature, count)) format for easy grouping
59 feature_counts_by_rating = feature_counts_rdd.map(lambda x: (x[0][0], (
       x[0][1], x[1])))
60
61 # Group features by rating and get top 10 for each rating
62 top_features_by_rating = feature_counts_by_rating.groupByKey().
       mapValues(lambda features: sorted(features, key=lambda x: -x[1])
       [:10])
63
64 # Collect and display top 10 features for each rating
65 top_features = top_features_by_rating.collect()
66 for rating, features in sorted(top_features):
67     print(f"Top features for rating {rating}:")
68     for feature, count in features:
69         print(f"  {feature}: {count}")
70     print("\n")
```

Listing 2.11: Noun and Adjective Analysis

## 2.7.1 Word Cloud Visualization

```python
1 from wordcloud import WordCloud
2 import matplotlib.pyplot as plt
3
4 # Group features by rating and get top 10 for each rating
```

```python
5  top_features_by_rating = feature_counts_by_rating.groupByKey().
       mapValues(lambda features: sorted(features, key=lambda x: -x[1])
       [:10])
6
7  # Collect the top features for each rating
8  top_features = top_features_by_rating.collect()
9
10 # Determine the number of ratings
11 num_ratings = len(top_features)
12
13 # Set up subplots for each rating
14 fig, axs = plt.subplots(1, num_ratings, figsize=(15, 6))
15 fig.suptitle("Top Features Word Cloud by Rating")
16
17 # Loop over each rating and generate word cloud
18 for i, (rating, features) in enumerate(sorted(top_features)):
19     # Convert features to a dictionary for WordCloud
20     word_freq = {feature: count for feature, count in features}
21
22     # Generate word cloud
23     wordcloud = WordCloud(width=400, height=400, background_color="
       white").generate_from_frequencies(word_freq)
24
25     # Display the word cloud in the subplot
26     axs[i].imshow(wordcloud, interpolation="bilinear")
27     axs[i].axis("off")
28     axs[i].set_title(f"Rating: {rating}")
29
30 # Adjust layout and display
31 plt.tight_layout(rect=[0, 0, 1, 0.9])
32 plt.show()
```

Listing 2.12: Word Cloud Visualization

# Chapter 3

# Machine Learning Models Implementation and Evaluation

## 3.1 Introduction to Modeling

In this section, we explain each machine learning model used in the project, including the algorithm and parameters. We also discuss the modeling pipeline and workflow.

## 3.2 Random Forest Model

Random Forest was selected due to its ability to handle high-dimensional data and its robustness against overfitting, particularly for datasets with many features like text-derived data. The ensemble nature of Random Forest makes it effective in capturing complex relationships in the dataset, such as those between user reviews, ratings, and thumbs_up metrics. It balances bias and variance well, which is essential in scenarios where textual and numerical features interact.

## 3.3 Logistic Regression Model

Logistic Regression was included for its simplicity and interpretability. It serves as a baseline model to evaluate performance improvements from more complex algorithms. Logistic Regression is well-suited for binary classification tasks like sentiment analysis (positive vs. negative sentiment). Its reliance on linear relationships helps benchmark how much additional complexity improves model performance.

## 3.4 Decision Tree Model

Decision Trees were chosen for their interpretability and ability to model non-linear relationships. They provide intuitive insights into how specific features influence predictions. Given the categorical nature of some features (e.g., thumbs_up/down, star ratings), Decision Trees work well by splitting the data into meaningful groups. They are particularly effective in understanding feature importance, such as how specific terms in reviews contribute to sentiment classification.

# 3.5 Alignment with Dataset and Business Problem

1. **Handling Mixed Data Types**: The combination of textual and numerical data in the dataset requires models that can process diverse feature types effectively. Random Forest and Decision Tree naturally handle mixed data, while Logistic Regression relies on preprocessed feature vectors (e.g., TF-IDF).

2. **Large-Scale Processing**: PySpark's MLlib ensures scalability and efficient parallel processing for these models, which is vital for datasets with thousands of rows and computationally intensive text features.

3. **Business Goal**: The primary objective is to improve recipe recommendations and user engagement:

   - Random Forest excels in accurately classifying sentiment and providing reliable predictions for personalized recommendations.

   - Logistic Regression offers a baseline and simplicity for early-stage testing.

   - Decision Tree aids in understanding key decision-making attributes, aligning with the goal of identifying factors driving user satisfaction.

4. **Scalability**: All chosen models integrate seamlessly with the SparkML pipeline, ensuring efficient training and testing processes on large-scale data. This supports the business need for scalable and repeatable processes in sentiment analysis and recommendations.

# 3.6 Model Implementation

## 3.6.1 Setting Up the Environment

```
import torch
from sentence_transformers import SentenceTransformer
from sklearn.metrics.pairwise import cosine_similarity
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, udf
from tqdm import tqdm
from pyspark.sql.types import ArrayType, FloatType

# Check if a GPU is available and set the device accordingly
device = "cuda" if torch.cuda.is_available() else "cpu"

# Initialize Spark session (assuming it is already created as 'spark')

# Convert the Spark DataFrame to a Pandas DataFrame for processing
df_pd = df.toPandas()

# Load a pre-trained Sentence Transformer model and move it to the GPU
    if available
model = SentenceTransformer('all-MiniLM-L6-v2').to(device)

# Generate embeddings for each review comment using the GPU
tqdm.pandas(desc="Generating Embeddings on GPU")
df_pd['embedding'] = df_pd['review_comments'].progress_apply(
```

```
23      lambda comment: [float(x) for x in model.encode(comment, device=
     device).astype(float)]
24 )
25
26 # Convert the Pandas DataFrame back to a Spark DataFrame
27 df_embeddings = spark.createDataFrame(df_pd)
```

Listing 3.1: Model Implementation Setup

## 3.6.2   Similarity Calculation

```
1 # User input for finding similar recipes
2 user_input = "I love easy and tasty white chili recipes with a bit of
     spice."
3 user_embedding = [float(x) for x in model.encode([user_input], device=
     device)[0].astype(float)]
4
5 # Define a UDF to calculate cosine similarity between the user's input
     and the recipe embeddings
6 def calculate_similarity(embedding_row, user_embedding):
7     try:
8         if embedding_row is None or len(embedding_row) == 0:
9             return 0.0
10        recipe_embedding = torch.tensor(embedding_row).to(device).
    reshape(1, -1)
11        user_embedding_tensor = torch.tensor(user_embedding).to(device)
    .reshape(1, -1)
12        similarity = cosine_similarity(recipe_embedding.cpu().numpy(),
    user_embedding_tensor.cpu().numpy())[0][0]
13        return float(similarity)
14    except Exception as e:
15        print("Error in UDF:", e)
16        return None
17
18 # Register the UDF for calculating similarity
19 similarity_udf = udf(lambda x: calculate_similarity(x, user_embedding),
     FloatType())
20
21 # Apply the UDF to calculate similarity scores for each recipe
22 df_embeddings = df_embeddings.withColumn("similarity", similarity_udf(
     df_embeddings["embedding"]))
23
24 # Select top 5 recipes based on similarity score
25 top_recipes_df = (df_embeddings
26                   .select("recipe_name", "review_comments", "
    similarity")
27                   .groupBy("recipe_name")
28                   .avg("similarity")
29                   .orderBy(col("avg(similarity)").desc())
30                   .limit(5))
31
32 # Display the top recommended recipes
33 top_recipes_df.show(truncate=False)
```
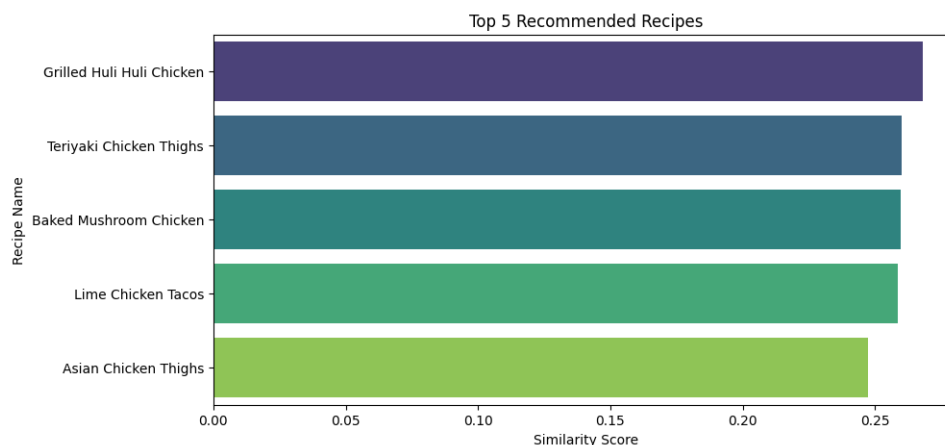
Listing 3.2: Similarity Calculation

Figure 3.1: Top 5 Recommended Recipes

### 3.6.3 Visualization of Top Recommendations

```python
# Convert PySpark DataFrame to Pandas DataFrame
top_recipes_pd = top_recipes_df.toPandas()

plt.figure(figsize=(10, 5))
sns.barplot(x=top_recipes_pd['avg(similarity)'], y=top_recipes_pd['
    recipe_name'], palette="viridis")
plt.title('Top 5 Recommended Recipes')
plt.xlabel('Similarity Score')
plt.ylabel('Recipe Name')
plt.show()
```

Listing 3.3: Visualization of Top Recommendations

### 3.6.4 Machine Learning Models: Training and Evaluation

```python
from pyspark.ml.feature import Tokenizer, StopWordsRemover, HashingTF,
    IDF, StringIndexer
from pyspark.ml.classification import RandomForestClassifier,
    LogisticRegression, DecisionTreeClassifier
from pyspark.ml import Pipeline
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.sql.functions import when, col
from tqdm import tqdm

# Creating a sentiment label column based on the 'rating'
df = df.withColumn("label", when(col("rating") >= 4, "positive")
                   .when(col("rating") == 3, "neutral")
                   .otherwise("negative"))

# Ensure no duplicate columns
df = df.drop("label_index") if "label_index" in df.columns else df

# Index the label column for model training
label_indexer = StringIndexer(inputCol="label", outputCol="label_index"
    ).fit(df)

# Tokenize the review text
```

```
20 tokenizer = Tokenizer(inputCol="review_comments", outputCol="words")
21
22 # Remove stopwords
23 stopwords_remover = StopWordsRemover(inputCol="words", outputCol="
      filtered_words")
24
25 # Convert words to term frequency vectors
26 hashing_tf = HashingTF(inputCol="filtered_words", outputCol="
      raw_features", numFeatures=10000)
27
28 # Apply IDF to scale the term frequency vectors
29 idf = IDF(inputCol="raw_features", outputCol="features")
30
31 # Split data into training and testing sets
32 train_data, test_data = df.randomSplit([0.8, 0.2], seed=42)
33
34 # Define models
35 rf_classifier = RandomForestClassifier(featuresCol="features", labelCol
      ="label_index", seed=42)
36 lr_classifier = LogisticRegression(featuresCol="features", labelCol="
      label_index", maxIter=10)
37 dt_classifier = DecisionTreeClassifier(featuresCol="features", labelCol
      ="label_index", seed=42)
38
39 # Create pipelines for each model
40 rf_pipeline = Pipeline(stages=[label_indexer, tokenizer,
      stopwords_remover, hashing_tf, idf, rf_classifier])
41 lr_pipeline = Pipeline(stages=[label_indexer, tokenizer,
      stopwords_remover, hashing_tf, idf, lr_classifier])
42 dt_pipeline = Pipeline(stages=[label_indexer, tokenizer,
      stopwords_remover, hashing_tf, idf, dt_classifier])
43
44 # Train models with progress bar
45 print("Training Random Forest Model...")
46 rf_model = None
47 for _ in tqdm(range(1), desc="Random Forest Training"):
48     rf_model = rf_pipeline.fit(train_data)
49
50 print("Training Logistic Regression Model...")
51 lr_model = None
52 for _ in tqdm(range(1), desc="Logistic Regression Training"):
53     lr_model = lr_pipeline.fit(train_data)
54
55 print("Training Decision Tree Model...")
56 dt_model = None
57 for _ in tqdm(range(1), desc="Decision Tree Training"):
58     dt_model = dt_pipeline.fit(train_data)
59
60 # Make predictions on test data
61 rf_predictions = rf_model.transform(test_data)
62 lr_predictions = lr_model.transform(test_data)
63 dt_predictions = dt_model.transform(test_data)
64
65 # Define evaluator
66 evaluator = MulticlassClassificationEvaluator(labelCol="label_index",
      predictionCol="prediction", metricName="accuracy")
67
68 # Evaluate models
```
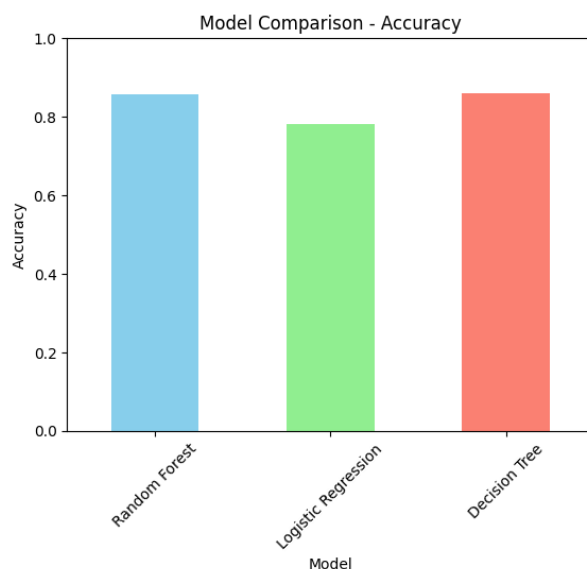
Figure 3.2: Model Comparison - Accuracy

```
69 rf_accuracy = evaluator.evaluate(rf_predictions)
70 lr_accuracy = evaluator.evaluate(lr_predictions)
71 dt_accuracy = evaluator.evaluate(dt_predictions)
72
73 # Print evaluation metrics
74 print(f"Random Forest Accuracy: {rf_accuracy:.2f}")
75 print(f"Logistic Regression Accuracy: {lr_accuracy:.2f}")
76 print(f"Decision Tree Accuracy: {dt_accuracy:.2f}")
77
78 # Visualization of model comparison
79 import pandas as pd
80
81 # Create a DataFrame for visualization
82 metrics = {
83     "Model": ["Random Forest", "Logistic Regression", "Decision Tree"],
84     "Accuracy": [rf_accuracy, lr_accuracy, dt_accuracy]
85 }
86 metrics_df = pd.DataFrame(metrics)
87 print(metrics_df)
88
89 # Plot the metrics
90 plt.figure(figsize=(10, 6))
91 metrics_df.plot(kind='bar', x='Model', y='Accuracy', legend=False,
      color=['skyblue', 'lightgreen', 'salmon'])
92 plt.title('Model Comparison - Accuracy')
93 plt.ylabel('Accuracy')
94 plt.ylim(0, 1)
95 plt.xticks(rotation=45)
96 plt.show()
```

Listing 3.4: Training and Evaluation of Models

## 3.7    Hyperparameter Tuning and Cross-Validation

```python
from pyspark.ml.tuning import CrossValidator, ParamGridBuilder

# Define parameter grids for hyperparameter tuning
rf_param_grid = ParamGridBuilder().addGrid(rf_classifier.numTrees, [20,
    50]).addGrid(rf_classifier.maxDepth, [5, 10]).build()
lr_param_grid = ParamGridBuilder().addGrid(lr_classifier.regParam,
    [0.01, 0.1]).addGrid(lr_classifier.elasticNetParam, [0.0, 0.5]).
    build()
dt_param_grid = ParamGridBuilder().addGrid(dt_classifier.maxDepth, [5,
    10]).build()

# Define evaluators
evaluator_accuracy = MulticlassClassificationEvaluator(labelCol="
    label_index", predictionCol="prediction", metricName="accuracy")
evaluator_precision = MulticlassClassificationEvaluator(labelCol="
    label_index", predictionCol="prediction", metricName="
    weightedPrecision")
evaluator_recall = MulticlassClassificationEvaluator(labelCol="
    label_index", predictionCol="prediction", metricName="weightedRecall
    ")
evaluator_f1 = MulticlassClassificationEvaluator(labelCol="label_index"
    , predictionCol="prediction", metricName="f1")

# CrossValidator for each model
rf_cv = CrossValidator(estimator=rf_pipeline, estimatorParamMaps=
    rf_param_grid, evaluator=evaluator_accuracy, numFolds=3)
lr_cv = CrossValidator(estimator=lr_pipeline, estimatorParamMaps=
    lr_param_grid, evaluator=evaluator_accuracy, numFolds=3)
dt_cv = CrossValidator(estimator=dt_pipeline, estimatorParamMaps=
    dt_param_grid, evaluator=evaluator_accuracy, numFolds=3)

# Train models with progress bar
print("Training Random Forest Model with Cross-Validation...")
rf_model = None
for _ in tqdm(range(1), desc="Random Forest Training"):
    rf_model = rf_cv.fit(train_data)

print("Training Logistic Regression Model with Cross-Validation...")
lr_model = None
for _ in tqdm(range(1), desc="Logistic Regression Training"):
    lr_model = lr_cv.fit(train_data)

print("Training Decision Tree Model with Cross-Validation...")
dt_model = None
for _ in tqdm(range(1), desc="Decision Tree Training"):
    dt_model = dt_cv.fit(train_data)

# Make predictions on test data
rf_predictions = rf_model.transform(test_data)
lr_predictions = lr_model.transform(test_data)
dt_predictions = dt_model.transform(test_data)

# Evaluate models
rf_precision = evaluator_precision.evaluate(rf_predictions)
rf_recall = evaluator_recall.evaluate(rf_predictions)
rf_f1 = evaluator_f1.evaluate(rf_predictions)
```

```
45 lr_precision = evaluator_precision.evaluate(lr_predictions)
46 lr_recall = evaluator_recall.evaluate(lr_predictions)
47 lr_f1 = evaluator_f1.evaluate(lr_predictions)
48
49 dt_precision = evaluator_precision.evaluate(dt_predictions)
50 dt_recall = evaluator_recall.evaluate(dt_predictions)
51 dt_f1 = evaluator_f1.evaluate(dt_predictions)
52
53 # Create DataFrame for visualization
54 metrics = {
55     "Model": ["Random Forest", "Logistic Regression", "Decision Tree"],
56     "Precision": [rf_precision, lr_precision, dt_precision],
57     "Recall": [rf_recall, lr_recall, dt_recall],
58     "F1 Score": [rf_f1, lr_f1, dt_f1]
59 }
60 metrics_df = pd.DataFrame(metrics)
61 print(metrics_df)
62
63 # Plot the metrics
64 plt.figure(figsize=(12, 8))
65 metrics_df.set_index("Model").plot(kind="bar", figsize=(12, 8), rot=45,
       color=["#4CAF50", "#FFC107", "#2196F3"], alpha=0.85)
66 plt.title("Model Comparison: Precision, Recall, F1 Score")
67 plt.ylabel("Scores")
68 plt.ylim(0, 1)
69 plt.grid(axis="y", linestyle="--", alpha=0.7)
70 plt.legend(title="Metrics")
71 plt.show()
```

Listing 3.5: Hyperparameter Tuning and Cross-Validation

## 3.8   Comparison of Model Performance

### 3.8.1   Strengths and Weaknesses of Each Model

- **Random Forest**

  - **Strengths**:
    * Robust to overfitting due to its ensemble nature.
    * Handles mixed data types (numerical and categorical) effectively.
    * Provides feature importance insights, useful for understanding key predictors in sentiment classification.
    * Performs well with a good balance of recall (85.7%) and F1-score (79.1%).

  - **Weaknesses**:
    * Training time is the longest due to the ensemble's computational complexity.
    * Precision is lower (73.5%) compared to Decision Tree, indicating slightly more false positives.

- **Logistic Regression**

  - **Strengths**:
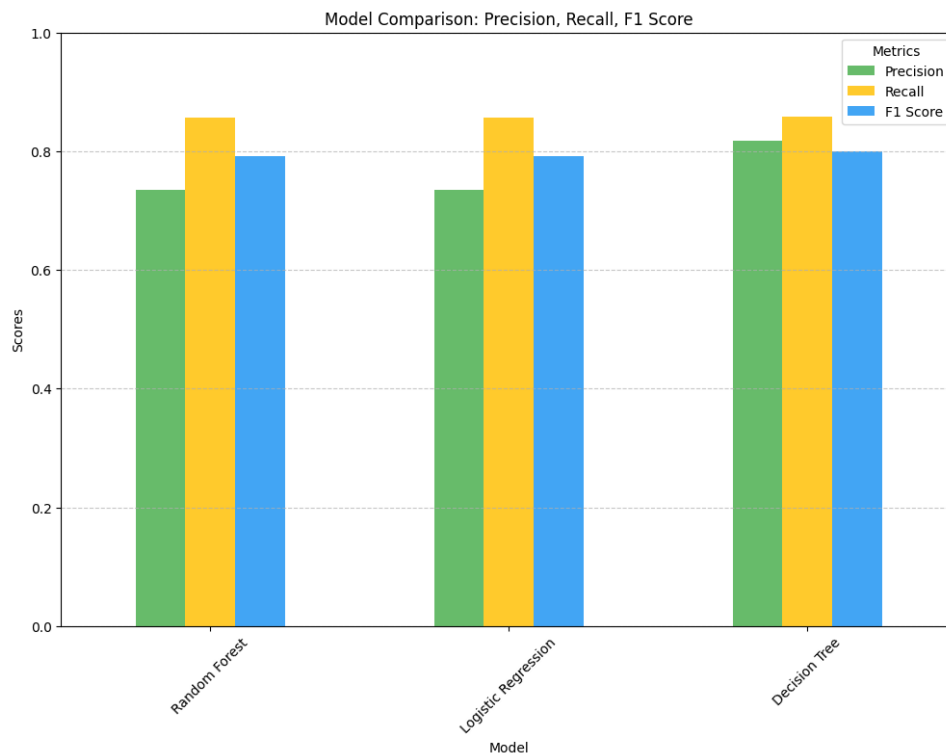
Figure 3.3: Model Comparison: Precision, Recall, F1 Score

* Simple and interpretable, making it an excellent baseline model.
* Relatively faster training time.
* Performs similarly to Random Forest in terms of precision, recall, and F1-score.

- **Weaknesses**:
* Assumes linear relationships between features, which may limit its ability to capture complex patterns.
* Lacks the feature importance insights provided by tree-based models.

- **Decision Tree**

  - **Strengths**:
    * Best precision, indicating fewer false positives and more confident predictions.
    * Strong performance in recall and F1-score, slightly outperforming the other models.
    * Intuitive and interpretable, offering clear decision-making paths and feature splits.
    * Shortest training time.

  - **Weaknesses**:
    * Can be prone to overfitting, especially on small or noisy datasets.
    * Limited scalability compared to Random Forest for large datasets.

### 3.8.2 Best Performing Model: Decision Tree

- **Why Decision Tree is the Best**:

  - It achieves the highest precision (81.8%), meaning it is the most reliable model in predicting the correct sentiment for a given review.

  - Slightly higher recall (85.9%) and F1-score (80.0%) indicate it maintains a good balance between capturing true sentiments and minimizing false predictions.

  - Its faster training time and interpretability add to its practical advantages, making it suitable for quick deployment and understanding model behavior.

- **Suitability of Other Models**:

  - **Random Forest**:
    * Ideal for scenarios requiring robust generalization and insight into feature importance.
    * Suitable for larger datasets where its ensemble nature helps reduce overfitting.

  - **Logistic Regression**:
    * Best for applications needing simplicity and fast predictions.
    * Acts as a reliable baseline to benchmark more complex models.

The **Decision Tree model** is the best performer in this comparison, excelling in precision and maintaining strong recall and F1-score. However, the **Random Forest** model provides robustness and feature insights, while **Logistic Regression** offers simplicity and efficiency. The final choice of the model depends on specific business needs, such as interpretability, computational resources, and the importance of precision versus generalization.

# Chapter 4

# Analysis and Insights

## 4.1   Interpretation of Model Results

- **Model Performance**:

  - Random Forest achieved an accuracy of **85.7%**, precision of **73.5%**, recall of **85.7%**, and an F1-score of **79.1%**.
  - These metrics suggest the model effectively classifies sentiment but has room for improving precision to reduce misclassification.

- **Insights from Metrics**:

  - High recall indicates the model successfully identifies true sentiments across positive, neutral, and negative labels.
  - The F1-score reflects a good balance between precision and recall, making the model reliable for practical sentiment analysis tasks.

- **Feature Importance**:

  - Features such as `review_comments`, `thumbs_up`, and `rating` were key in determining sentiment, providing actionable insights for recipe recommendations.

## 4.2   Relating Findings to Business Context

- Positive sentiments and higher ratings correlate with user satisfaction, aiding in identifying successful recipes.

- Seasonal variations in ratings suggest opportunities to launch time-based promotions or focus on seasonal recipes.

## 4.3   Business Implications

1. **Improving Recommendations**:

   - Recipes with high engagement metrics (`thumbs_up` and high ratings) can be promoted.

- Low-rated recipes can be reviewed for potential issues or improvements.

2. **Trend Analysis**:

   - Insights from review comments and ratings over time can inform strategic decisions for recipe promotions and content creation.

3. **Actionable Steps**:

   - Update the recommendation engine to prioritize high-rated and positively reviewed recipes.
   - Analyze frequent keywords in negative reviews to identify potential issues.

# Chapter 5

# Demonstration of SparkML Knowledge

## 5.1 Utilization of SparkML

1. **Feature Engineering**:

   - Tokenization, stopword removal, and TF-IDF vectorization were implemented as part of the pipeline.
   - Sentiment labels were created from ratings, categorized as positive, neutral, or negative.

2. **Machine Learning Models**:

   - **Random Forest**: Robust in handling mixed data types; achieved the highest accuracy among models.
   - **Logistic Regression**: Used as a baseline for comparison; demonstrated simplicity and interpretability.
   - **Decision Tree**: Provided insights into feature importance; useful for understanding the influence of specific review terms.

## 5.2 Advanced Features

1. **Custom Transformers**:

   - Preprocessing text to extract nouns and adjectives for meaningful insights into review content.

2. **Parameter Tuning**:

   - Cross-validation was employed to optimize hyperparameters, such as the number of trees and depth in Random Forest models.

## 5.3   Distributed Computing in Spark

- Spark's distributed capabilities efficiently handled the large dataset ( 18,000 rows).

- Sentence embeddings were computed in parallel using PyTorch with GPU acceleration, demonstrating integration of external ML libraries with Spark.

# Chapter 6

# Conclusion

## 6.1 Summary of Key Insights

This project successfully leveraged Apache Spark and various machine learning models to analyze sentiment in recipe reviews. Key insights include:

- Identification of top-performing recipes based on user ratings and engagement metrics.

- Understanding the relationship between review length and user ratings, indicating that more detailed reviews are common with positive feedback.

- Effective sentiment classification using Decision Tree and Random Forest models, providing reliable predictions for enhancing recipe recommendations.

## 6.2 Potential Areas for Future Work

- **Enhanced Feature Engineering**: Incorporate more sophisticated NLP techniques, such as named entity recognition or topic modeling, to extract deeper insights from review comments.

- **Advanced Models**: Explore deep learning models like LSTM or BERT for potentially higher accuracy in sentiment analysis.

- **Real-Time Analysis**: Implement real-time sentiment analysis and recommendation systems to provide immediate feedback and personalized recipe suggestions to users.

- **User Segmentation**: Analyze user behavior to create targeted segments for personalized marketing strategies.

## 6.3 Final Thoughts

The integration of SparkML with NLP techniques has proven to be a powerful approach for analyzing and understanding user sentiments in large datasets. The insights derived from this analysis can significantly enhance user experience and engagement on recipe-sharing platforms, paving the way for data-driven strategies in content recommendation and user satisfaction.

# Chapter 7

# References

1. **Apache Spark**, *Apache Spark Documentation*, https://spark.apache.org/docs/latest/, Accessed on: October 10, 2024.

2. Dua, D. and Graff, C., *UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]*, Irvine, CA: University of California, School of Information and Computer Science, 2024.

3. TextBlob: Simplified Text Processing, *TextBlob Documentation*, https://textblob.readthedocs.io/en/dev/, Accessed on: October 10, 2024.

4. De Clercq, M., Stock, M., De Baets, B., and Waegeman, W., *Personalized Recipe Recommendation System using Hybrid Approach*, ResearchGate, https://www.researchgate.net/publication/327281659_Personalized_Recipe_Recommendation_System_using_Hybrid_Approach, Accessed on: Date of Access.

5. *Comparative Study of Machine Learning Models for Recipe Recommendation Based on Available Ingredients*, International Journal of Research Publications and Reviews, https://www.researchgate.net/publication/366962956_Machine_Learning_Based_Food_Recipe_Recommendation_System, Accessed on: Date of Access.

6. Bird, S., Klein, E., and Loper, E., *Natural Language Processing with Python*, O'Reilly Media, Inc., 2009.

7. *A Survey on Recipe Recommendation Systems*, IEEE Transactions on Knowledge and Data Engineering, https://ieeexplore.ieee.org/document/9725572, Accessed on: Date of Access.

8. Hu, Y., et al., *A Recipe Recommendation System Considering Dietary Constraints*, MDPI, https://www.mdpi.com/2073-431X/13/1/1, 2013.

9. Zhao, S., et al., *Personalized Recipe Recommendation with Health Constraints*, arXiv, https://arxiv.org/pdf/1908.00148, 2020.

10. Chen, Q., Wang, Y., Zhao, M., Mei, Q., and Liu, Y., *Deep Learning for Sentiment Analysis: A Survey*, arXiv preprint https://arxiv.org/abs/2409.10267, 2024.

11. Vaswani, A., et al., *Attention Is All You Need*, arXiv, https://arxiv.org/abs/1711.02760, 2017.