

---

## **Implementation of Web Scraping and Data Mining for Performance Evaluation of PT Ceria Multimedia Services**

Ravel Yanuartha<sup>1\*</sup>, Dini Fakta Sari<sup>2</sup>

<sup>1,2</sup>Universitas Teknologi Digital Indonesia Yogyakarta, Fakultas Teknologi Informasi, Program Studi Informatika, Jl. Raya Janti Jl. Majapahit No.143, Jaranan, Banguntapan, Kec. Banguntapan, Kabupaten Bantul, Daerah Istimewa Yogyakarta, Indonesia

---

### **Keywords**

*Application; Information; Software; System*

### **\*Corresponding Author:**

[ravel.yanuartha@students.utdi.ac.id](mailto:ravel.yanuartha@students.utdi.ac.id)

### **Abstract**

This research aims to implement web scraping techniques to collect testimonial data from the CeriaMultimedia website and perform sentiment analysis to evaluate service quality. The collected data consists of a limited number of testimonials, which are then processed through text preprocessing stages including case folding, tokenizing, filtering, and stemming. The sentiment classification process is conducted using machine learning methods based on TF-IDF weighting and classification algorithms. Due to the limited dataset, the analysis results are used primarily to demonstrate the implementation process rather than to draw generalized conclusions. The results show that the sentiment categories obtained include positive, negative, and neutral sentiments, although not all categories consistently appear in the testing phase. This research highlights the effectiveness of web scraping and text processing techniques while also indicating the need for a larger dataset to improve evaluation accuracy in future studies.

---

## **1. Introduction**

Computer technology plays an important role in supporting information processing across various sectors, including business services [1]. In the service industry, evaluating service performance is essential to ensure customer satisfaction and continuous improvement. However, service performance assessment in many multimedia equipment rental businesses is still conducted manually and subjectively, which can lead to inefficiencies and biased evaluations. Along with the increasing number of multimedia equipment rentals, meetings, training, CAT tests, workshops, seminars, product launches, and gatherings at Ceria Multimedia, PT Ceria Multimedia operates a service system that supports customers who require multimedia equipment and meeting facilities [2]. However, collecting service-related information directly from customers requires a considerable amount of time, especially when comparisons with other multimedia service providers are needed. Therefore, an automated system is required to collect service information efficiently. One approach that can be used is web scraping, which enables the automatic retrieval of semi-structured data from websites, such as customer testimonials published on the Ceria Multimedia website [3].

Data mining is a process that involves extracting meaningful information from large datasets using statistical and machine learning techniques. The Knowledge Discovery in Databases (KDD) process supports the identification of useful patterns from complex data automatically and accurately [4]. Machine learning techniques further assist in analyzing textual data to support decision-making based on discovered patterns [5]. In the context of service evaluation, sentiment analysis is widely used to classify customer opinions into positive, negative, or neutral categories [6]. However, previous studies mainly focus on social media data and do not specifically address service performance evaluation using testimonials from official company websites.

Based on these considerations, this research aims to answer the following research question: *How can web scraping and data mining techniques be implemented to evaluate the service performance of PT Ceria Multimedia based on customer testimonials?* This study focuses on implementing web scraping to collect testimonial data from the official PT Ceria Multimedia website and applying data mining techniques to analyze the collected data for service performance evaluation. The scraping process is conducted using the Python programming language on the official Ceria Multimedia website. The results of this study are expected to provide insights into customer perceptions of service quality and support objective evaluation of service performance at PT Ceria Multimedia [7].

## **2. Research Method**

The methods used in this study were selected based on their suitability for processing textual testimonial data and evaluating service performance. Web scraping was chosen as it enables efficient and automated collection of customer testimonials from official websites, which have been widely used in previous sentiment analysis studies [1], [2]. Data mining techniques were applied to extract meaningful patterns from the collected data to support objective service performance evaluation.

In this research, the implementation of web scraping and data mining is conducted to evaluate the service performance of PT Ceria Multimedia. The data used in this study consist of customer testimonials obtained from the official website of PT Ceria Multimedia. It should be noted that the dataset used in this research is limited in size, consisting of only ten testimonial entries. Therefore, the results of this study are intended to demonstrate the implementation of the proposed method rather than to provide generalizable conclusions.

The research methodology consists of the following stages:

### **2.1 Data collection**

Data collection was carried out using web scraping techniques to automatically extract customer testimonials from the Ceria Multimedia website. The scraped data were stored in CSV format to facilitate further processing [7], [8]. Although web scraping allows efficient and automated data acquisition, the testimonial data used in this study were collected solely from the official website of PT Ceria Multimedia. As a result, the dataset represents only customers who voluntarily submitted testimonials and may not fully capture the perspectives of all service users. Consequently, this limitation should be considered when interpreting the evaluation results, as the collected data may not comprehensively reflect overall customer satisfaction.

This study considers ethical aspects in the data collection process using web scraping techniques. The data collected consist solely of publicly available customer testimonials published on the official website of PT Ceria Multimedia. No private, sensitive, or personally identifiable information was accessed or stored during the scraping process. The data were collected strictly for academic research purposes and were not used for commercial exploitation. Furthermore, the scraping process was conducted in a manner that did not interfere with the normal operation of the website, thereby respecting data ownership and responsible data usage principles.

### **2.2 Data Pre-Processing**

The pre-processing stage aims to prepare textual data for analysis. This stage includes case folding, text cleaning, tokenizing, stopword removal, and stemming to reduce noise and standardize the text data [9].

### 2.3 Data Transformation

Data transformation was performed using the Term Frequency–Inverse Document Frequency (TF-IDF) method to convert textual data into numerical feature vectors. TF-IDF assigns weights to words based on their importance in the document collection, calculated, with the formula [10]:

$$TF - IDF(t, d) = TF(t, d) \times IDF(t) \quad (1)$$

$$TF(t, d) = \frac{\text{jumlah kemunculan kata } t \text{ dalam dokumen } d}{\text{jumlah total kata dalam dokumen } d} \quad (2)$$

$$IDF(t) = \log \frac{N}{n_t} \quad (3)$$

where:

N = total number of documents

$n_t$  = number of documents containing word t

### 2.4 Core Process

Core Process: in this core process, performance evaluation (average score, satisfaction) uses the Weighted Average Scoring method, with the formula [1], [11]:

$$\text{average score} = \frac{\sum_{i=1}^n R_i}{n} \quad (4)$$

where:

$R_i$  = rating of the i-th customer

n = total number of testimonials

### 2.5 Result Evaluation

The evaluation stage applies sentiment analysis to classify customer testimonials into positive, negative, and neutral categories. The classification results are used to observe the proportion of customer satisfaction toward the services provided by PT Ceria Multimedia [4].

### 2.6 Visualization

The final results are presented through visualizations in the form of charts to clearly display sentiment distribution and service evaluation outcomes [10]. Due to the limited dataset size, advanced validation techniques such as k-fold cross-validation and the use of larger datasets are suggested for future research to improve model reliability and evaluation accuracy.

## 3. Result and Discussions

This study only uses testimonial data collected from the official website of PT Ceria Multimedia. In future research, data can be collected from multiple sources such as social media platforms or online review services to obtain more diverse and representative customer opinions.

### 3.1 Data Collection

Data Collection is the process of gathering information or data from specific sources to be used for research, analysis, or model development. In the context of your research or project, data collection refers to extracting testimonial data from the PT Ceria Multimedia website so that it can be further analyzed (for example, to evaluate service performance).

Table 1. Testimonial Description

No	Name	Testimonial Description	Rating Score
1	Rosa Mega Ariani	Peralatan yang dimiliki sangat lengkap dan jumlahnya banyak. Sangat puas dengan layanan yang baik dan ramah. Semoga Ceria Multimedia tetap eksis menjadi mitra selamanya.	5
2	Dicka Herdiansyah	Dengan kantor cabang tersebar di berbagai kota di Indonesia, sangat membantu kami para owner Even Organizer. Dapat job event tinggal call aja. Cepet tanpa ribet!!	5
3	Irma Septiani Putri	Alhamdulillah seluruh rangkaian acara berjalan lancar. Terima kasih Ceria Multimedia, next event bisa bekerjasama lagi.	5
4	I Made Krisna	Profesional,, kelihatan dari caranya menangani event sangat berpengalaman. Sy rekomendasikan buat teman-teman semua. Bravo Ceria Multimedia..	5
5	Andi Yosapea	Tim Ceria Multimedia solid banget. Semua memiliki etos kerja tinggi. Lembur sampai larut malam mereka jalani demi suksesnya acara. Trims ya udah full support event kami.	5
6	Jhonatan Rai Sambian	Pelayanan sangat bagus dan memuaskan terima kasih ceria multimedia	5
7	Femmy Arista Caniago	Sangat ramah sekali.	5
8	Endri Nur Rohman	Keren sekali event kami berjalan dengan lancar, namun ada beberapa yang perlu di perbaiki adalah cara komunikasi yang baik dan benar.	4
9	Niken Kania	Ceria Multimedia peralatannya lengkap banget, tapi ada beberapa alat yang harus di update sesuai dengan perkembangan jaman dan supaya Ceria Multimedia makin maju.	3
10	Harry Joel	Team ceriamultimedia solid, namun ada beberapa yang perlu di perbaiki supaya teamnya makin solid.	3

### 3.2 Pre-Processing

The preprocessing stage is used to process existing data to avoid interference from inconsistent data. The goal is to achieve high accuracy in the classification output, and this stage transforms disorganized data into well-organized data. Text preprocessing involves six processes: case folding, cleaning, tokenizing, stop word removal, and steaming.

#### 3.2.1 Case Folding

The first preprocessing step in this study is case folding. This case folding step converts the characters of words in the dataset from uppercase to lowercase. After the coding process is performed, the resulting dataset is no longer capitalized. Table 2 below shows the results after the text was converted using case folding:

Table 2. Case Folding

No	Case Folding
1	peralatan yang dimiliki sangat lengkap dan jumlahnya banyak. sangat puas dengan layanan yang baik dan ramah. semoga ceria multimedia tetap eksis menjadi mitra selamanya.
2	dengan kantor cabang tersebar di berbagai kota di indonesia, sangat membantu kami para owner even organizer. dapat job event tinggal call aja. cepet tanpa ribet!!
3	alhamdulillah seluruh rangkaian acara berjalan lancar. terima kasih ceria multimedia, next event bisa bekerjasama lagi.
4	profesional,, kelihatan dari caranya menangani event sangat berpengalaman. sy rekomendasikan buat teman-teman semua. bravo ceria multimedia..

No	Case Folding
5	tim ceria multimedia solid banget. semua memiliki etos kerja tinggi. lembur sampai larut malam mereka jalani demi suksesnya acara. trims ya udah full support event kami.
6	pelayanan sangat bagus dan memuaskan terima kasih ceria multimedia
7	sangat ramah sekali.
8	keren sekali event kami berjalan dengan lancar, namun ada beberapa yang perlu di perbaiki adalah cara komunikasi yang baik dan benar.
9	ceria multimedia peralatannya lengkap banget, tapi ada beberapa alat yang harus di update sesuai dengan perkembangan jaman dan supaya ceria multimedia makin maju.
10	team ceriamultimedia solid, namun ada beberapa yang perlu di perbaiki supaya teamnya makin solid.

### 3.2.2 Cleaning

Next, we perform the cleaning phase using Python to remove data attributes. This includes removing unnecessary characters from the text, such as HTML tags and emoticons mixed with comments within the text, which can interfere with the data analysis process. Table 3 shows the results after removing the mixed characters in the comments, as follows:

Table 3. Cleaning

No	Cleaning
1	peralatan yang dimiliki sangat lengkap dan jumlahnya banyak sangat puas dengan layanan yang baik dan ramah semoga ceria multimedia tetap eksis menjadi mitra selamanya
2	dengan kantor cabang tersebar di berbagai kota di indonesia sangat membantu kami para owner even organizer dapat job event tinggal call aja cepet tanpa ribet
3	alhamdulillah seluruh rangkaian acara berjalan lancar terima kasih ceria multimedia next event bisa bekerjasama lagi
4	profesional kelihatan dari caranya menangani event sangat berpengalaman sy rekomendasikan buat teman semua bravo ceria multimedia
5	tim ceria multimedia solid banget semua memiliki etos kerja tinggi lembur sampai larut malam mereka jalani demi suksesnya acara trims ya udah full support event kami
6	pelayanan sangat bagus dan memuaskan terima kasih ceria multimedia
7	sangat ramah sekali
8	keren sekali event kami berjalan dengan lancar namun ada beberapa yang perlu di perbaiki adalah cara komunikasi yang baik dan benar
9	ceria multimedia peralatannya lengkap banget tapi ada beberapa alat yang harus di update sesuai dengan perkembangan jaman dan supaya ceria multimedia makin maju
10	team ceriamultimedia solid namun ada beberapa yang perlu di perbaiki supaya teamnya makin solid

### 3.2.3 Tokenizing

The next step is tokenization, where a sentence from a Google Maps comment or dataset is broken down into word chunks, or tokens. This is done to identify word occurrences and prepare the text for word-level analysis. Table 4 shows the results after tokenization:

Table 4. Tokenizing

No	Tokenizing
1	['peralatan', 'yang', 'dimiliki', 'sangat', 'lengkap', 'dan', 'jumlahnya', 'banyak', 'sangat', 'puas', 'dengan', 'layanan', 'yang', 'baik', 'dan', 'ramah', 'semoga', 'ceria', 'multimedia', 'tetap', 'eksis', 'menjadi', 'mitra', 'selamanya']

No	Tokenizing
2	['dengan', 'kantor', 'cabang', 'tersebar', 'di', 'berbagai', 'kota', 'di', 'indonesia', 'sangat', 'membantu', 'kami', 'para', 'owner', 'even', 'organizer', 'dapat', 'job', 'event', 'tinggal', 'call', 'aja', 'cepat', 'tanpa', 'ribet']
3	['alhamdulillah', 'seluruh', 'rangkaian', 'acara', 'berjalan', 'lancar', 'terima', 'kasih', 'ceria', 'multimedia', 'next', 'event', 'bisa', 'bekerjasama', 'lagi']
4	['profesional', 'kelihatan', 'dari', 'caranya', 'menangani', 'event', 'sangat', 'berpengalaman', 'sy', 'rekomendasikan', 'buat', 'temanteman', 'semua', 'bravo', 'ceria', 'multimedia']
5	['tim', 'ceria', 'multimedia', 'solid', 'banget', 'semua', 'memiliki', 'etos', 'kerja', 'tinggi', 'lembur', 'sampai', 'larut', 'malam', 'mereka', 'jalani', 'demi', 'suksesnya', 'acara', 'trims', 'ya', 'udah', 'full', 'support', 'event', 'kami']
6	['pelayanan', 'sangat', 'bagus', 'dan', 'memuaskan', 'terima', 'kasih', 'ceria', 'multimedia']
7	['sangat', 'ramah', 'sekali']
8	['keren', 'sekali', 'event', 'kami', 'berjalan', 'dengan', 'lancar', 'namun', 'ada', 'beberapa', 'yang', 'perlu', 'di', 'perbaiki', 'adalah', 'cara', 'komunikasi', 'yang', 'baik', 'dan', 'benar']
9	['ceria', 'multimedia', 'peralatannya', 'lengkap', 'banget', 'tapi', 'ada', 'beberapa', 'alat', 'yang', 'harus', 'di', 'update', 'sesuai', 'dengan', 'perkembangan', 'jaman', 'dan', 'supaya', 'ceria', 'multimedia', 'makin', 'maju']
10	['team', 'ceriamultimedia', 'solid', 'namun', 'ada', 'beberapa', 'yang', 'perlu', 'di', 'perbaiki', 'supaya', 'timnya', 'makin', 'solid']

### 3.2.4 Stopword

Stopword removal is the process of removing stopwords (common words that appear frequently but provide little value in text analysis) from documents or text data. Its main purpose is to simplify text data and focus analysis on more relevant and informative words. Words such as: yang, dan, di, ke, dari, ini, itu, pada, adalah, dengan, sebagai. Table 5 shows the results after stopwords removal:

Table 5. Stopword

No	Stopword
1	['peralatan', 'dimiliki', 'sangat', 'lengkap', 'jumlahnya', 'banyak', 'sangat', 'puas', 'layanan', 'baik', 'ramah', 'semoga', 'ceria', 'multimedia', 'tetap', 'eksis', 'menjadi', 'mitra', 'selamanya']
2	['kantor', 'cabang', 'tersebar', 'berbagai', 'kota', 'indonesia', 'sangat', 'membantu', 'kami', 'para', 'owner', 'even', 'organizer', 'dapat', 'job', 'event', 'tinggal', 'call', 'aja', 'cepat', 'tanpa', 'ribet']
3	['alhamdulillah', 'seluruh', 'rangkaian', 'acara', 'berjalan', 'lancar', 'terima', 'kasih', 'ceria', 'multimedia', 'next', 'event', 'bisa', 'bekerjasama', 'lagi']
4	['profesional', 'kelihatan', 'caranya', 'menangani', 'event', 'sangat', 'berpengalaman', 'sy', 'rekomendasikan', 'buat', 'temanteman', 'semua', 'bravo', 'ceria', 'multimedia']
5	['tim', 'ceria', 'multimedia', 'solid', 'banget', 'semua', 'memiliki', 'etos', 'kerja', 'tinggi', 'lembur', 'sampai', 'larut', 'malam', 'mereka', 'jalani', 'demi', 'suksesnya', 'acara', 'trims', 'ya', 'udah', 'full', 'support', 'event', 'kami']
6	['pelayanan', 'sangat', 'bagus', 'memuaskan', 'terima', 'kasih', 'ceria', 'multimedia']
7	['sangat', 'ramah', 'sekali']
8	['keren', 'sekali', 'event', 'kami', 'berjalan', 'lancar', 'namun', 'ada', 'beberapa', 'perlu', 'perbaiki', 'cara', 'komunikasi', 'baik', 'benar']
9	['ceria', 'multimedia', 'peralatannya', 'lengkap', 'banget', 'tapi', 'ada', 'beberapa', 'alat', 'harus', 'update', 'sesuai', 'perkembangan', 'jaman', 'supaya', 'ceria', 'multimedia', 'makin', 'maju']
10	['team', 'ceriamultimedia', 'solid', 'namun', 'ada', 'beberapa', 'perlu', 'perbaiki', 'supaya', 'timnya', 'makin', 'solid']

### 3.2.5 Stemming

Stemming is the process of changing an affixed word into a base word by removing prefixes, suffixes, insertions and confixes contained in the word. The main goal of stemming is to simplify variations of words that have the

same meaning so that it can increase the effectiveness of text analysis and reduce the number of word features in text data processing:

Table 6. Stemming

No	Stemming
1	['alat', 'milik', 'sangat', 'lengkap', 'jumlah', 'banyak', 'sangat', 'puas', 'layan', 'baik', 'ramah', 'moga', 'ceria', 'multimedia', 'tetap', 'eks', 'jadi', 'mitra', 'lama']
2	['kantor', 'cabang', 'sebar', 'bagai', 'kota', 'indonesia', 'sangat', 'bantu', 'kami', 'para', 'owner', 'even', 'organizer', 'dapat', 'job', 'event', 'tinggal', 'call', 'aja', 'cepat', 'tanpa', 'ribet']
3	['alhamdulillah', 'seluruh', 'rangkai', 'acara', 'jalan', 'lancar', 'terima', 'kasih', 'ceria', 'multimedia', 'next', 'event', 'bisa', 'bekerjasama', 'lagi']
4	['profesional', 'lihat', 'cara', 'tangan', 'event', 'sangat', 'alam', 'sy', 'rekomendasi', 'buat', 'temanteman', 'semua', 'bravo', 'ceria', 'multimedia']
5	['tim', 'ceria', 'multimedia', 'solid', 'banget', 'semua', 'milik', 'etos', 'kerja', 'tinggi', 'lembur', 'sampai', 'larut', 'malam', 'mereka', 'jalan', 'demi', 'sukses', 'acara', 'trims', 'ya', 'udah', 'full', 'support', 'event', 'kami']
6	['layan', 'sangat', 'bagus', 'muas', 'terima', 'kasih', 'ceria', 'multimedia']
7	['sangat', 'ramah', 'sekali']
8	['keren', 'sekali', 'event', 'kami', 'jalan', 'lancar', 'namun', 'ada', 'beberapa', 'perlu', 'baik', 'cara', 'komunikasi', 'baik', 'benar']
9	['ceria', 'multimedia', 'alat', 'lengkap', 'banget', 'tapi', 'ada', 'beberapa', 'alat', 'harus', 'update', 'sesuai', 'kembang', 'jaman', 'supaya', 'ceria', 'multimedia', 'makin', 'maju']
10	['team', 'ceriamultimedia', 'solid', 'namun', 'ada', 'beberapa', 'perlu', 'baik', 'supaya', 'timnya', 'makin', 'solid']

### 3.3 Sentiment Labeling

Sentiment labeling was conducted to assign sentiment categories to each customer testimonial. The labeling process classified testimonials into three sentiment classes: positive, negative, and neutral. This study employed a lexicon-based approach, where sentiment labels were determined based on the presence of sentiment-bearing words in the text.

Positive sentiment is indicated by words such as “ramah”, “puas”, and “baik”, which reflect customer satisfaction with the services provided. Negative sentiment is identified through words expressing dissatisfaction or complaints, while neutral sentiment represents testimonials that contain descriptive information without clear emotional polarity. Given the limited size of the dataset, sentiment labeling was performed to ensure consistency and clarity in sentiment classification, serving as the basis for subsequent analysis and evaluation.

### 3.4 Visualization

This visualization process was carried out to understand the words or terms that generated the highest frequency in the research data. shows the results, which include three sentiments: positive, neutral, and negative. The larger the word cloud, the higher its frequency.



program uses the words "tapi,,namun, perlu" so that the words that often appear are the words ""tapi,,namun, perlu " as shown in the image above.

### 3.5 Word Weighting (TF-IDF)

Word weighting using the Term Frequency–Inverse Document Frequency (TF-IDF) method was applied to identify important words in testimonial data collected from the Ceria Multimedia website. The dataset used in this study consists of a limited number of testimonial records obtained through web scraping. Each column in the TF-IDF matrix represents unique terms found in the testimonial corpus after the preprocessing stage. TF-IDF assigns higher weights to terms that frequently appear in a specific testimonial but occur less frequently across the entire dataset, indicating their importance in representing customer opinions [10]. Based on the TF-IDF results, the top weighted terms were further analyzed according to sentiment categories (positive, neutral, and negative).

- For sentiment positive , high-weighted terms such as "*ramah, puas, baik*" indicate that customers generally appreciate service quality, staff professionalism, and punctuality provided by PT Ceria Multimedia. These terms reflect positive customer perceptions related to service reliability and overall satisfaction.
- In the sentiment neutral category, terms like "*ceriamultimedia, acara, layanan* " tend to have moderate TF-IDF weights. These words suggest that customers perceive the service as acceptable but not exceptional, indicating opportunities for further service improvement.
- Meanwhile, sentiment negative is characterized by high-weighted terms such as "*tapi,namun, perlu*". These terms reveal customer dissatisfaction related to service responsiveness and technical performance, highlighting aspects that require attention and improvement by the service provider [4].

Overall, the analysis of top weighted terms per sentiment using TF-IDF provides insights into customer perceptions by identifying key factors influencing satisfaction and dissatisfaction. This information can support PT Ceria Multimedia in evaluating service performance and prioritizing improvements based on customer feedback [12].

```
0,Rosa Mega Ariani,alat yang milik sangat lengkap dan jumlah banyak sangat puas dengan layan yang baik dan ramah moga ceria multimedia tetap eks jadi mitra lama,4.477355348654375
1,Dicka Herdiansyah,dengan kantor cabang sebar di bagai kota di indonesia sangat bantu kami para owner even organizer dapat job event tinggal call aja cepet tanpa ribet,4.83688774963430
1
2,Irma Septiani Putri,alhamdulillah seluruh rangkai acara jalan lancar terima kasih ceria multimedia next event bisa bekerjasama lagi,3.7990038890408266
3,I Made Krisna,profesional lihat dari cara tangan event sangat alan sy rekomendasi buat teman semua bravo ceria multimedia,3.9135395207787935
4,Andi Yosapea,tim ceria multimedia solid banget semua milik etos kerja tinggi lembur sampai larut malam mereka jalan demi sukses acara trims ya udah full support event kami,5.0303901910837086
5,Jhonatan Rai Sambian,layan sangat bagus dan muas terima kasih ceria multimedia,2.924134524272697
6,Femmy Arista Caniago,sangat ramah sekali,1.7108314911361913
7,Endri Nur Rohman,keren sekali event kami jalan dengan lancar namun ada beberapa yang perlu di baik adalah cara komunikasi yang baik dan benar,4.2262678111878484
8,Niken Kania,ceria multimedia alat lengkap banget tapi ada beberapa alat yang harus di update sesuai dengan kembang jaman dan supaya ceria multimedia makin maju,4.33927416163892
9,Harry Joel,team ceriamultimedia solid namun ada beberapa yang perlu di baik supaya teamnya makin solid,3.46460111084397
```

Figure 4. Word Weighting Results (TF-IDF)

Shows that the first line with the text "tim ceria multimedia solid banget semua milik etos kerja tinggi lembur sampai larut malam mereka jalan demi sukses acara trims ya udah full support event kami " has a total TF-IDF value of 5.0303901910837086, which is one of the highest values among the examples shown.

### 3.6 Core Process

The core process is the main stage of this research, which aims to evaluate the service performance of PT Ceria Multimedia based on the results of sentiment analysis derived from customer testimonials. The data used at this stage consist of sentiment classification results that have undergone text preprocessing and word weighting using the TF-IDF method, ensuring that the data are suitable for quantitative service performance evaluation. Service performance evaluation is conducted using the Weighted Average Scoring method, where each sentiment category (positive, neutral, and negative) is assigned a specific weight according to its

contribution to customer satisfaction. This approach produces a single evaluation score that represents the overall service performance based on customer perceptions.

Table 7. Sentiment Weight

Sentiment	R value
Positive	3
Negative	2
Neutral	1

Table 8. Sentiment Interpretation

R Value Range	Interpretation
2.34 – 3.00	Good / Satisfactory
1.67 – 2.33	Moderate
1.00 – 1.66	Poor

Table 9. Sentiment Count

Sentiment	Count
Positive	6
Negative	2
Neutral	2

Total weighted score:

$$(6 \times 3) + (2 \times 2) + (2 \times 1) = 24 \quad (5)$$

Total testimonials:

$$6 + 2 + 2 = 10 \quad (6)$$

Weighted average score:

$$Score = \frac{24}{10} = 2.40 \quad (7)$$

Based on Table 8, the weighted average score of 2.40 falls within the range 2.34–3.00, which indicates that the service performance of PT Ceria Multimedia is categorized as Good / Satisfactory. This result shows that positive sentiment dominates customer testimonials, reflecting generally favorable customer perceptions of the services provided.

### 3.7 Evaluation

In this study, a total of 10 data samples were used and divided into training and testing sets using an 80:20 ratio, where eight samples were used for training and two samples for testing. Considering the limited size of the dataset, k-fold cross-validation was applied to improve the reliability of the evaluation by allowing each data sample to be used as both training and testing data across multiple folds. This approach helps reduce bias caused by a single data split and provides a more robust estimation of model performance.

Model performance was evaluated using accuracy, precision, recall, and F1-score, which provide a more comprehensive assessment than accuracy alone. The evaluation results are illustrated in Figure 7. The relatively low performance obtained in this study can be attributed to the small dataset size, which limits the model's ability to learn diverse linguistic patterns and may cause bias toward certain sentiment classes. In addition, the imbalance of sentiment categories in the dataset affects the classifier's ability to correctly predict minority classes.

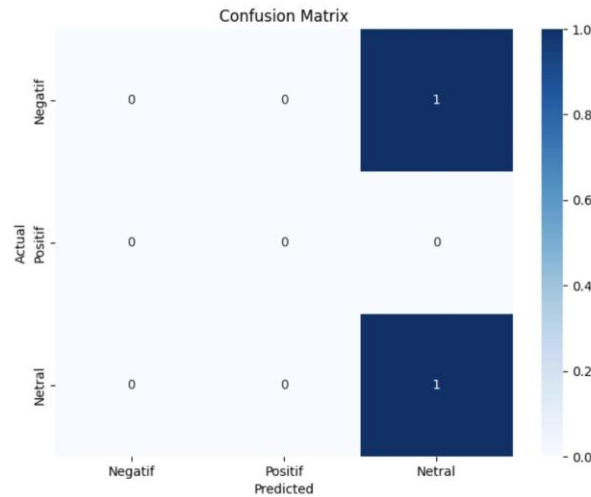


Figure 5. Confusion Matrix

```

Classification Report:
              precision    recall  f1-score   support

   Negatif      0.00      0.00      0.00         1
   Netral       0.50      1.00      0.67         1

 accuracy      0.50
 macro avg     0.25      0.50      0.33
 weighted avg  0.25      0.50      0.33

Manual Weighted Averages:
Precision: 0.25
Recall: 0.50
F1-Score: 0.33
Total Support: 2.0
    
```

Figure 6. Result Accuracy

The accuracy value is calculated as follows:

$$Accuracy = \frac{(0.0 \times 1) + (1.00 \times 1)}{1 + 1} = \frac{1.00}{2} = 0.50 \quad (8)$$

Figure 6 shows the application of the Naïve Bayes algorithm in this study using a ratio of 80:20, resulting in an accuracy of 0.50%.

### 3.8 Discussion

This study implements web scraping and data mining techniques to evaluate customer service performance at PT Ceria Multimedia through sentiment analysis of customer testimonials. The applied methodology follows common practices in sentiment analysis research, particularly in the stages of text preprocessing, feature weighting, and result interpretation. Text preprocessing plays an important role in sentiment analysis

performance. In this research, preprocessing stages include case folding, cleaning, tokenizing, stopword removal, and stemming. Previous studies have shown that these preprocessing steps significantly influence sentiment classification results and reduce noise in textual data [9]. Similar preprocessing pipelines were also applied in sentiment analysis studies using Naive Bayes classifiers [1].

TF-IDF was employed as a word weighting method to represent textual features and identify important terms in customer testimonials. Several studies have demonstrated that TF-IDF is effective in highlighting dominant terms that contribute to sentiment polarity[12], [13], [14]. In this research, TF-IDF results reveal key terms frequently associated with customer satisfaction and dissatisfaction. Positive sentiment testimonials tend to emphasize terms related to service quality and equipment readiness, while negative sentiment highlights issues related to service delays or unmet expectations. These findings indicate that TF-IDF can effectively capture customer perceptions toward service performance.

Compared to previous sentiment analysis studies that utilized large datasets collected from social media platforms such as Twitter and Instagram [2], [3], [6], this study is limited by the relatively small number of testimonial records obtained from a single company website. The limited dataset size may affect the robustness and generalizability of the sentiment distribution results. Similar limitations were acknowledged in other studies with constrained datasets, where sentiment proportions were influenced by data imbalance[14], [15]. Despite these limitations, the results of this study are consistent with prior research indicating that sentiment analysis can provide meaningful insights into customer satisfaction when combined with appropriate preprocessing and feature weighting techniques [16], [17], [18]. The visualization of sentiment proportions helps PT Ceria Multimedia understand overall customer perception and identify areas requiring service improvement.

Future research is recommended to collect data from multiple sources and increase dataset size to improve classification reliability. Additionally, the use of k-fold cross-validation and comparison with other machine learning classifiers such as Support Vector Machine may enhance evaluation accuracy, as demonstrated in previous studies[4], [13], [19].

However, the limited size of the dataset and its reliance on a single data source may affect the robustness of the analysis results. With only ten testimonial records, the sentiment distribution and performance evaluation may be sensitive to individual data points and may not generalize well to broader customer populations. Consequently, the findings of this study should be interpreted as indicative rather than conclusive, and future studies are encouraged to utilize larger and more diverse datasets to improve model reliability and evaluation accuracy.

#### **4. Conclusions and Future Works**

Based on the research results, it can be concluded that the web scraping process was successfully implemented to retrieve testimonial data from the CeriaMultimedia website and store it in a structured format. The obtained data can then be processed through text preprocessing stages, such as data cleaning, stemming, and word weighting using the TF-IDF method. The sentiment analysis performed was able to group testimonials into positive, negative, and neutral sentiment categories. However, due to the relatively small amount of data used, the model evaluation results did not fully represent all sentiment categories consistently in the test data. Therefore, the results of this study emphasize the implementation method and the sentiment analysis process flow rather than the model's accuracy level.

#### **5. References**

- [1] A. K. Qorita and F. Rahma, "Analisis Sentimen Berdasarkan Aspek pada Tempat Wisata di Daerah Istimewa Yogyakarta," *AUTOMATA*, vol. 3, no. 1, 2022, [Online]. Available: <https://journal.uii.ac.id/AUTOMATA/article/view/21906>

- [2] P. Pandunata, C. K. Ananta, and Y. Nurdiansyah, "Analisis Sentimen Opini Publik Terhadap Pekan Olahraga Nasional Pada Instagram Menggunakan Metode Naïve Bayes Classifier," *INFORMAL Informatics Journal*, vol. 7, no. 2, pp. 146–156, 2022, doi: 10.19184/isj.v7i2.33928.
- [3] R. Sulastiyono, A. Setiawan, and S. Nugroho, "Sentimen Analisis Pembatalan Indonesia Menjadi Tuan Rumah Piala Dunia U-20 Menggunakan Metode Naïve Bayes," *Journal of Information System Research*, vol. 4, no. 4, pp. 1387–1394, 2023, doi: 10.47065/josh.v4i4.3737.
- [4] F. P. Herlambang and D. Avianto, "Analisis Sentimen Opini Pengguna Twitter Terhadap Tragedi Kanjuruhan Malang dengan Metode Support Vector Machine," *Jurnal Media Informatika Budidarma*, vol. 7, no. 4, pp. 1727–1739, 2023, doi: 10.30865/mib.v7i4.6332.
- [5] M. Kholilullah, M. Martanto, and U. Hayati, "Analisis Sentimen Pengguna Twitter (X) Tentang Piala Dunia Usia 17 Menggunakan Metode Naive Bayes," *JATI*, vol. 8, no. 1, pp. 392–398, 2024, doi: 10.36040/jati.v8i1.8378.
- [6] R. Rasiban and S. Riyadi, "Analisis Sentimen Opini Masyarakat Terhadap Stadion Jakarta International Stadium (JIS) Pada Twitter Dengan Perbandingan Metode Naive Bayes dan Support Vector Machine," *Jurnal Sains dan Teknologi*, vol. 5, no. 3, pp. 1010–1017, 2024, doi: 10.55338/saintek.v5i3.2790.
- [7] C. A. Cholik, "Perkembangan Teknologi Informasi Komunikasi (ICT) dalam Berbagai Bidang," *Jurnal Fakultas Teknik UNISA Kuningan*, vol. 2, no. 2, pp. 39–46, 2021.
- [8] A. Purwansyah, A. Afriyudi, and S. Suyanto, "Perancangan dan Implementasi Sistem Informasi Pelaporan Masyarakat untuk Kerusakan Jalan di Palembang Menggunakan Google Maps API," *Jurnal Nasional Ilmu Komputer*, vol. 1, no. 4, pp. 175–182, 2020, doi: 10.47747/jurnalnik.v1i4.164.
- [9] S. Khairunnisa, "Pengaruh Text Preprocessing terhadap Analisis Sentimen Komentar Twitter," *Jurnal Media Informatika Budidarma*, vol. 5, no. 2, pp. 406–414, 2021.
- [10] K. M. Mahendra, M. H. Murdiansyah, D. T. Lhaksana, "Analisis Sentimen Tweet COVID-19 Menggunakan K-Nearest Neighbors dengan TF-IDF dan CountVectorizer," *DIKE: Jurnal Ilmu Multidisiplin*, vol. 1, no. 2, pp. 37–43, 2023.
- [11] A. Pebdika, R. Herdiana, and D. Solihudin, "Klasifikasi Menggunakan Metode Naive Bayes untuk Menentukan Calon Penerima PIP," *JATI*, vol. 7, no. 1, pp. 452–458, 2023.
- [12] A. H. Ayatullah, "Analisis Sentimen Penilaian Masyarakat terhadap Pelayanan Rumah Sakit Muhammadiyah Lamongan Menggunakan TF-IDF dan Naive Bayes," *Jurnal Informatika Medis*, vol. 2, no. 1, pp. 27–33, 2024, doi: 10.52060/im.v2i1.2198.
- [13] K. M. Mahendra, M. H. Murdiansyah, and D. T. Lhaksana, "Analisis Sentimen Tweet COVID-19 Menggunakan K-Nearest Neighbors dengan TF-IDF dan CountVectorizer," *DIKE: Jurnal Ilmu Multidisiplin*, vol. 1, no. 2, pp. 37–43, 2023, doi: 10.69688/dike.v1i2.35.
- [14] T. I. Alfawas, "Penerapan TF-IDF untuk Analisis Sentimen Ulasan Game Bus Simulator Indonesia," *Innovative Journal of Social Science Research*, vol. 4, no. 5, pp. 3177–3193, 2024.
- [15] A. Munawaroh, "Sentiment Analysis dengan Naive Bayes terhadap Risiko Pembangunan IKN," *JATI*, 2024, doi: 10.36040/jati.v8i1.8454.
- [16] E. Suryati, Styawati, and A. A. Aldino, "Analisis Sentimen Transportasi Online Menggunakan Word2Vec dan SVM," *Jurnal Teknologi dan Sistem Informasi*, vol. 4, no. 1, pp. 96–106, 2023.
- [17] N. P. Husain, "Analisis Sentimen Ulasan TikTok Berbasis TF-IDF dan SVM," *JSCE*, vol. 5, no. 1, pp. 91–102, 2024, doi: 10.61628/jsce.v5i1.1105.

- [18] H. C. Husada and A. S. Paramita, "Analisis Sentimen Maskapai Penerbangan di Twitter Menggunakan SVM," *Teknika*, vol. 10, no. 1, pp. 18–26, 2021, doi: 10.34148/teknika.v10i1.311.
- [19] H. Apriyani and K. Kurniati, "Perbandingan Metode Naive Bayes dan SVM dalam Klasifikasi Penyakit Diabetes," *Jurnal Information Technology Ampera*, vol. 1, no. 3, pp. 133–143, 2020, doi: 10.51519/journalita.volume1.issue3.year2020.page133-143.