

Performance Evaluation of Sentiment Classification Models: A Comparative Study of NBC, SVM, and DT with SMOTE

Yerik Afrianto Singgalen*

Faculty of Business Administration and Communication, Tourism Department, Atma Jaya Catholic University of Indonesia, Jakarta, Indonesia

Email: yerik.afrianto@atmajaya.ac.id

Email Penulis Korespondensi: yerik.afrianto@atmajaya.ac.id

Abstract—This research explores the performance of sentiment classification models, namely Naive Bayes Classifier (NBC), Decision Tree (DT), and Support Vector Machine (SVM), using the CRISP-DM methodology in the context of digital content analysis and data mining. The analysis was conducted on a SMOTE dataset in Rapidminer, yielding significant performance metrics. The NBC model achieved an accuracy of 86.98% \pm 0.96%, precision of 100.00% \pm 0.00%, recall of 78.82% \pm 1.55%, and f-measure of 88.15% \pm 0.97%, with an AUC of 0.657 \pm 0.203. Similarly, the DT model exhibited an accuracy of 93.20% \pm 0.42%, precision of 90.87% \pm 0.64%, recall of 98.88% \pm 0.31%, and f-measure of 94.70% \pm 0.31%, with an AUC of 0.918 \pm 0.006. Furthermore, the SVM model demonstrated an accuracy of 96.80% \pm 0.65%, precision of 98.99% \pm 0.28%, recall of 95.77% \pm 1.03%, and f-measure of 97.35% \pm 0.55%, with an AUC of 0.994. These findings highlight the efficacy of these models in accurately classifying sentiments within digital content, suggesting their suitability for various data mining applications. Recommendations for future research include exploring ensemble methods, continuous model updating, alternative sampling techniques, feature engineering approaches, and collaboration with domain experts to enhance real-world applicability.

Keywords: DT; NBC; Sentiment Classification; SMOTE; SVM

1. INTRODUCTION

The advancement of digital media has significantly facilitated content creators in documenting the exploration outcomes of indigenous settlements and the livelihoods of traditional communities. Through various digital platforms and tools, creators intricately capture and portray indigenous peoples' rich cultural heritage and daily practices [1]. This digital revolution empowers creators to produce immersive visual and auditory content, enhancing the accessibility and preservation of indigenous knowledge and customs [2]–[6]. Moreover, the ease of dissemination afforded by digital media enables a broader audience to engage with and appreciate the intricate tapestry of indigenous lifestyles and societal dynamics [7]–[10]. As a result, the synergy between digital technology and cultural documentation fosters greater awareness, appreciation, and preservation of indigenous traditions for future generations.

The dissemination of documentary videos depicting the livelihoods and settlements of diverse communities across various online platforms incites viewers to engage with factors of livelihood strategies. These multimedia productions serve as potent tools for raising awareness and fostering understanding of the socio-economic dynamics inherent in different cultural contexts [11]–[14]. By intricately showcasing the daily lives and survival strategies of communities worldwide, these documentaries prompt viewers to contemplate the interconnectedness of environmental, social, and economic factors influencing livelihood choices [15]–[18]. Consequently, such multimedia publications contribute significantly to sustainable development and cultural preservation discourse.

This research investigates viewer responses by conducting sentiment analysis on review data posted in comment sections, elucidating the knowledge transfer processes facilitated by digital media and viewer behaviors in receiving digital information [19]–[24]. By scrutinizing sentiment patterns within user-generated content, such as comments and reviews, insights are gleaned into the efficacy of digital media in disseminating information and shaping audience perceptions [25]–[29]. Through this analytical lens, the intricate dynamics of knowledge dissemination and reception in the digital realm are comprehensively explored, shedding light on the evolving landscape of online communication and information consumption [30]–[33].

The primary aim of this study is to analyze viewer sentiment through review data on video content published by Kacong Explorer under the ID 9DIICbM76Bw. By focusing on the comments and feedback generated by viewers in response to Kacong Explorer's videos, the research seeks to discern prevailing sentiments and attitudes towards the content. Through systematic sentiment analysis, patterns in viewer perceptions and reactions were identified, providing valuable insights into audience engagement and the effectiveness of content delivery strategies employed by Kacong Explorer. Ultimately, this investigation contributes to a deeper understanding of audience preferences and behaviors within digital media consumption, informing future content creation and dissemination efforts.

The methodology employed in the data processing and analysis process is CRISP-DM (Cross-Industry Standard Process for Data Mining). This structured approach encompasses several phases: business understanding, data understanding, data preparation, modeling, evaluation, and deployment [34]. This research systematically navigates through each stage by adhering to the CRISP-DM framework, ensuring comprehensive coverage of the data mining process [35]. This methodological rigor facilitates efficient data processing and analysis, leading to robust insights and informed decision-making [36]. Consequently, using CRISP-DM enhances the reliability and validity of research findings, culminating in more impactful outcomes and actionable recommendations.

The urgency of this research lies in comprehending viewer sentiment regarding documentary videos showcasing explorations of indigenous villages. By delving into the perspectives and reactions of viewers towards such content, critical insights were gleaned into the efficacy of cultural representation and the resonance of indigenous narratives within contemporary digital media landscapes [37], [38]. This understanding is paramount for fostering cultural appreciation, preserving traditional knowledge, and promoting sustainable development initiatives that honor and respect indigenous communities [39], [40]. Consequently, the timely investigation of viewer sentiment holds significant implications for advancing discourse on cultural heritage preservation and fostering inclusive digital media representations.

This research's theoretical and practical implications are significant in advancing our understanding of digital media engagement and cultural representation. The findings contribute theoretically by enriching scholarly discourse on audience reception theories and digital media effects, particularly within the context of indigenous cultural preservation. Furthermore, the insights garnered from this study hold practical value for content creators, policymakers, and cultural advocates seeking to enhance the dissemination and reception of indigenous narratives in digital spaces. By elucidating viewer sentiments and preferences, this research informs strategies for crafting culturally sensitive and engaging content, fostering greater inclusivity and cultural appreciation in the digital realm. Ultimately, this research's theoretical and practical contributions enrich academic scholarship and practical initiatives to promote cultural diversity and preservation in the digital age.

The limitation of this research lies in its narrow focus on video dataset analysis and its reliance on the CRISP-DM methodology. While video data provides valuable insights into viewer sentiment and engagement, it inherently limits the scope of inquiry to visual media, potentially overlooking other digital content and audience interactions. Additionally, while CRISP-DM offers a structured framework for data mining and analysis, its applicability may be constrained by the complexity and diversity of data sources encountered in digital media research. Consequently, while this study offers valuable contributions within its defined parameters, future research endeavors should consider broadening the scope of analysis and adopting more flexible methodologies to comprehensively explore the multifaceted dynamics of digital media and audience reception.

2. RESEARCH METHODOLOGY

2.1 Gap Analysis

Gap analysis is a crucial stage in determining the extent to which research development on sentiment analysis contextualizes discussions around video content and processed review data. This research systematically identifies gaps and deficiencies within the current literature body through this analytical process, informing future research directions and methodological advancements. By assessing the adequacy of existing methodologies and theoretical frameworks, gap analysis catalyzes refining research approaches and enhancing the overall comprehensiveness and relevance of sentiment analysis studies in the digital media landscape. Consequently, the systematic undertaking of gap analysis facilitates a more nuanced understanding of the complexities inherent in analyzing sentiment within the context of video content and review data, thereby advancing scholarly discourse and practical applications in this field.

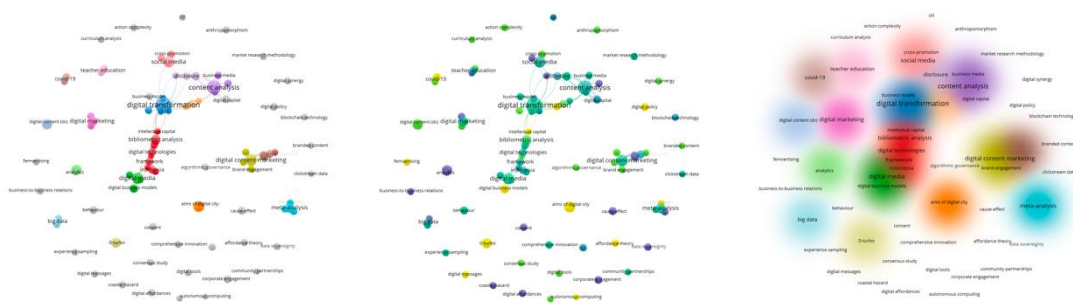


Figure 1. Gap Analysis using Vosviewer

Figure 1 shows the network, density, and overlay visualization using VosViewer. Based on the findings of the gap analysis, there is a clear need to enhance research on sentiment analysis for documentary videos to observe viewer behavior in comprehending thematically educational content, particularly within the context of this study's focus on the livelihood and settlements of the indigenous community of Kampung Naga. This research gains more profound insights into viewer engagement and knowledge absorption processes by directing attention toward the thematic understanding of educational content embedded within documentary videos. This emphasis on thematic analysis enriches our understanding of audience reception. It underscores the significance of contextualizing educational narratives within indigenous communities' cultural and socio-economic milieu like Kampung Naga. Consequently, advancing research in this direction holds the potential to foster more meaningful and impactful digital media representations of indigenous cultures while contributing to broader educational objectives.

2.2 Cross-Industry Standard Process for Data Mining (CRISP-DM)

The CRISP-DM framework comprises stages including business understanding, data understanding, modeling, evaluation, and deployment, tailored to the context of video content and its accompanying review data. This structured approach provides a systematic roadmap for conducting data mining and analysis, from gaining insights into the business objectives and requirements related to the video content to understanding the characteristics and quality of the available data. Subsequently, modeling techniques are applied to develop algorithms or models for sentiment analysis based on the data understanding phase. The evaluation stage assesses the effectiveness and accuracy of the models, while the deployment phase focuses on implementing the findings into practical applications or decision-making processes. Thus, the adaptability of the CRISP-DM framework enables this research to navigate the complexities of sentiment analysis within video content and systematically review data.

2.2.1 Business Understanding

During the business understanding stage, the video content and review data under analysis are specifically identified as the video with the ID 9DIICbM76Bw, published by Kacong Explorer on June 25, 2023. This video has garnered significant traction, accumulating 11,109,635 views and eliciting 13,414 comments. Such meticulous delineation of the video's attributes, including its publication date, view count, and comment volume, is pivotal for contextualizing subsequent stages of analysis within the CRISP-DM framework. By establishing a comprehensive understanding of the business objectives and the data landscape associated with the target video, this research effectively tailors the analytical approach to derive meaningful insights into viewer sentiment and engagement dynamics. Thus, the meticulous assessment of these foundational elements lays the groundwork for informed decision-making and strategic planning in subsequent phases of the analysis process.

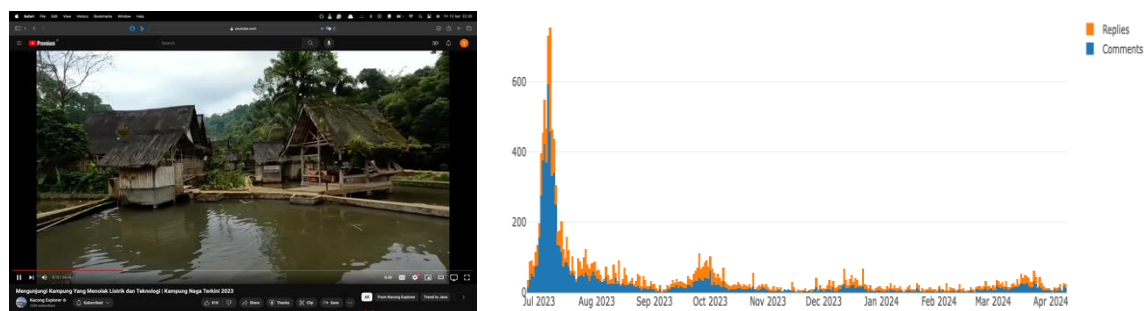


Figure 2. Post-Per-Day Statistic

Figure 2 shows the video and post-per-day statistics. Based on the post-per-day statistics, it is evident that the video received the highest number of comments on July 6, 2023, totaling 596 comments. The July 7 and July 4 dates are closely behind, with 459 and 423 comments, respectively. These statistics provide valuable insights into the temporal distribution of viewer engagement with the video content, highlighting specific dates when audience interaction peaked. Such detailed analysis of comment frequency across different dates enables this research to discern patterns and trends in viewer behavior. This facilitates a more nuanced understanding of audience dynamics and preferences within the digital media landscape. Consequently, leveraging this temporal granularity enhances the efficacy of sentiment analysis efforts and informs strategic decision-making processes to optimize audience engagement strategies and content delivery mechanisms.

As per the top ten poster data, it is apparent that @KacongExplorer emerges as the most prolific contributor with 317 posts, signifying a significant engagement and activity from the primary content creator. Conversely, the remaining posters exhibit substantially lower posting frequencies, with @nunungnuraeni8035, @user-id3yh8nh2v, @brahmaardiansyah3187, @sahrbinauk1349, @usepsuhendi60, @ZyoRoyz-yv5mf, @rusidahrusidah5656, @ekkiimanuel-iz4sm, and @adigunawan1243 contributing only a fraction of the total posts. Such data underscores the dominance of @KacongExplorer in shaping the discourse surrounding video content, highlighting the central role of content creators in driving audience engagement and interaction within digital media platforms. Consequently, this observation accentuates the significance of strategic collaborations and partnerships with influential content creators to optimize audience reach and engagement outcomes.



Figure 3. Top Ten Poster (Communaltyic)

Figure 3 shows the top ten posters of the dataset. Based on the top ten poster data, it becomes apparent that the highest number of comments originates from KaongExplorer, indicating active engagement from the content creator as the proprietor and manager of the channel in response to various viewer comments on the video. This trend underscores the proactive involvement of content creators in fostering dialogue and interaction within digital media platforms, cultivating a sense of community, and enhancing viewer satisfaction. Moreover, the substantial contribution of KaongExplorer to the comment section reflects a commitment to audience engagement and responsiveness, reinforcing the importance of creator-viewer interactions in shaping the overall viewer experience and perception of the content. Consequently, this observation highlights the pivotal role of content creators in fostering meaningful engagement and driving audience retention within digital media ecosystems.

2.2.2 Data Understanding

During the data understanding stage, the total downloaded review data amounted to 13,414 comments; however, following data cleansing procedures, the processed dataset was reduced to 8,939 comments. This reduction underscores the importance of data preprocessing techniques in refining and preparing the dataset for subsequent analysis. This research ensures the analytical outcomes' reliability and validity by eliminating irrelevant or redundant information and addressing data quality issues such as noise and inconsistencies. Consequently, this meticulous data cleansing process enhances the effectiveness of subsequent data mining and sentiment analysis efforts, facilitating more accurate and meaningful insights into viewer sentiments and engagement dynamics within the digital media landscape.

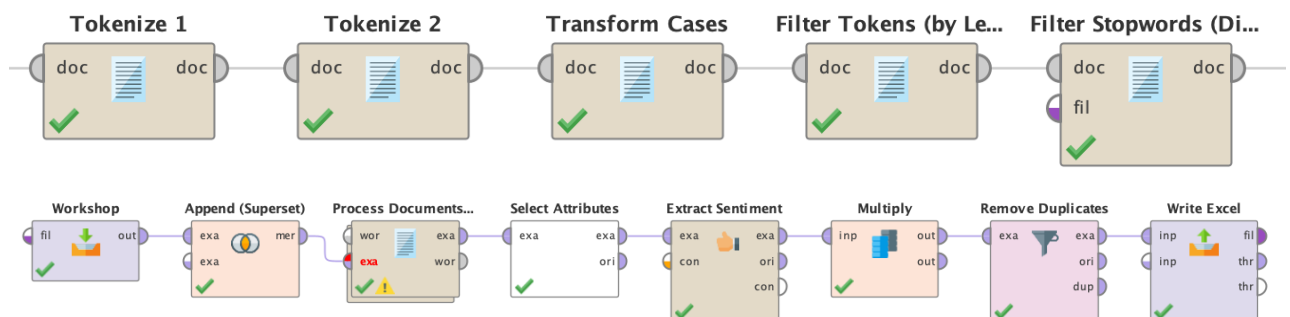


Figure 4. Data Cleaning and Extract Sentiment

Figure 4 shows the data cleaning process in Rapidminer using operator tokenize, transforms cases, filter tokens, and filter stopwords. Based on the sentiment extraction results, it is discernible that 158 instances were classified under the negative class out of the analyzed dataset, while the majority, totaling 8,782 review data, were categorized as belonging to the positive class. This disparity in sentiment distribution underscores the prevalence of positive sentiments among viewers towards the video content under scrutiny, indicating a favorable reception and engagement with the material. Such insights from sentiment analysis provide valuable feedback on audience perceptions and reactions and inform content creators and platform managers to refine the strategies to optimize viewer satisfaction and retention. Consequently, this nuanced understanding of sentiment dynamics enhances decision-making processes to enhance content relevance and effectiveness within the digital media landscape.

2.2.3 Modeling

During the modeling stage, the labeled data is subjected to further processing through performance testing of NBC, DT, and SVM algorithms. This phase involves applying machine learning techniques to train and evaluate predictive models based on the labeled dataset. By systematically testing the performance of different algorithms, this research assesses the effectiveness in accurately classifying sentiment within the dataset and identifying the most suitable approach for sentiment analysis. Consequently, this rigorous evaluation of algorithm performance enhances the robustness and reliability of the sentiment analysis process, ensuring that the derived insights accurately reflect the sentiment dynamics present within the analyzed data.

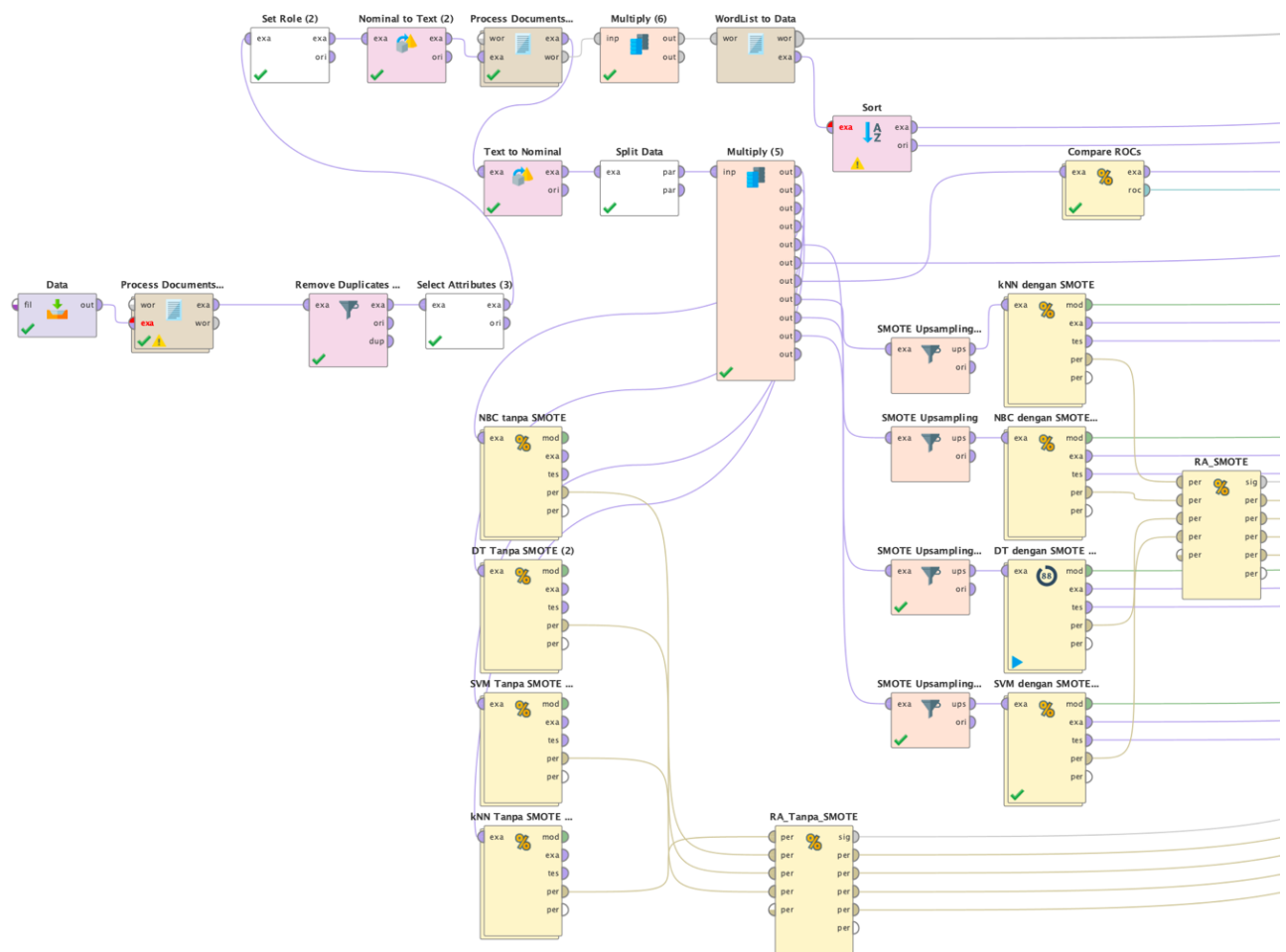


Figure 6. Implementation of NBC, DT, and SVM Models in Rapidminer

Figure 6 shows the modeling process in Rapidminer. In the modeling process, NBC, DT, and SMOTE are utilized, with SMOTE as the operator to address data imbalance. This technique is crucial for mitigating the impact of skewed class distributions commonly encountered in sentiment analysis tasks. Moreover, the data is partitioned into training and testing sets, with a split ratio of 70% for testing data and 30% for training data. This systematic approach ensures the adequacy of training data for model development while enabling robust evaluation of model performance on unseen test data. Consequently, by leveraging algorithmic techniques and rigorous data partitioning strategies, this research enhances the reliability and generalizability of sentiment analysis models, thereby facilitating more accurate and meaningful insights into viewer sentiments within digital media contexts.

2.2.4 Evaluation

During the evaluation stage, the best-performing algorithm or model is determined based on accuracy, precision, recall, F-measure, and area under the receiver operating characteristic curve (AUC). These evaluation metrics provide comprehensive insights into the performance of each model in accurately classifying sentiment within the dataset. By considering multiple performance indicators, this research assesses each model's effectiveness and reliability in capturing the nuances of viewer sentiments. Consequently, this systematic evaluation process enables this research to identify the most suitable algorithm or model for sentiment analysis, ensuring robust and accurate insights into viewer sentiments within digital media contexts.

2.2.5 Deployment

During the deployment stage, the model exhibiting the best performance is recommended as the most relevant for processing the dataset. However, within the contextual framework of documentary videos featuring the indigenous community of Kampung Naga, the dataset possesses distinctive characteristics that prompt viewers to engage with the content. Despite selecting the best-performing model for sentiment analysis, it is essential to acknowledge the unique contextual factors surrounding the dataset, such as cultural sensitivities and community-specific nuances, which may influence viewer interactions and feedback. Consequently, while deploying the sentiment analysis model, it is imperative to consider these contextual elements to ensure the relevance and accuracy of the analytical outcomes within the specific domain of indigenous cultural representation and exploration.

reviewers. Following closely are the 🙏 emoji with 20 occurrences, suggesting a notable presence of expressions of gratitude or prayer within the discussions. Additionally, positive emotions such as 😊, 😄, and 😁, each appearing 19, 18, and 18 times respectively, reflect a prevalent sentiment of happiness and amusement among the reviewers. Conversely, emojis like 😞 (16 occurrences) and 😞 (3 occurrences) indicate sadness or distress in the reviews. Furthermore, the presence of 👍 (13 occurrences) signifies approval or agreement, while 🤝 (11 occurrences) denotes mutual understanding or agreement among the reviewers. This comprehensive analysis of emoji usage offers valuable insights into viewers' nuanced emotional responses and sentiments towards the Kampung Naga documentary, underscoring the importance of visual cues in understanding audience engagement and perception within digital media contexts.

The words cloud and emojis, which prominently display positive expressions, symbolize the viewers' appreciation of the video content. This visual representation of frequently used positive words and emoticons underscores the prevailing sentiment of approval, satisfaction, and enjoyment experienced by the audience. Furthermore, the abundance of positive expressions within the words cloud and emoji usage reflects a collective endorsement and admiration for the thematic richness, storytelling, and cultural exploration depicted in the video. Consequently, the convergence of positive linguistic and visual cues highlights the effectiveness of the content in eliciting favorable responses and fostering audience engagement within digital media platforms.

KacongExplorer delineates the content of the video by emphasizing the adherence to intrinsic rules passed down through generations, aimed at preserving the environment and ensuring the continuity of ancestral traditions and culture amidst modernization pressures. This description underscores the intrinsic value placed on environmental conservation and cultural heritage within the indigenous community of Kampung Naga. The video highlights the community's resilience in safeguarding its identity and values by portraying the sustained commitment to ancestral practices despite modern influences. Consequently, this portrayal educates and informs viewers about the significance of cultural preservation and environmental stewardship and fosters a deeper appreciation for the harmonious coexistence of tradition and modernity within indigenous societies.

Furthermore, the evaluation results of the classification model with the best performance highlight the necessity of employing the Synthetic Minority Over-sampling Technique (SMOTE) to address imbalances within the dataset. This observation underscores the critical role of data preprocessing techniques in enhancing classification models' efficacy, particularly in skewed class distributions. By synthetically generating minority class samples, SMOTE mitigates the adverse effects of class imbalance, thereby accurately improving the model's ability to classify instances from underrepresented classes. Consequently, utilizing SMOTE emerges as a crucial strategy for enhancing the robustness and reliability of classification models in handling imbalanced datasets.

NBC with SMOTE		
PerformanceVector: accuracy: 86.98% +/- 0.96% (micro average: 86.98%) ConfusionMatrix: True: Negative Positive Negative: 3853 1302 Positive: 0 4845 AUC (optimistic): 1.000 +/- 0.000 (micro average: 1.000) (positive class: Positive) AUC: 0.657 +/- 0.203 (micro average: 0.657) (positive class: Positive) AUC (pessimistic): 0.789 +/- 0.015 (micro average: 0.789) (positive class: Positive) precision: 100.00% +/- 0.00% (micro average: 100.00%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3853 1302 Positive: 0 4845 recall: 78.82% +/- 1.55% (micro average: 78.82%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3853 1302 Positive: 0 4845 f_measure: 88.15% +/- 0.97% (micro average: 88.16%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3853 1302 Positive: 0 4845		
DT with SMOTE		
PerformanceVector: accuracy: 93.20% +/- 0.42% (micro average: 93.20%) ConfusionMatrix: True: Negative Positive Negative: 3242 69 Positive: 611 6078 AUC (optimistic): 0.997 +/- 0.002 (micro average: 0.997) (positive class: Positive) AUC: 0.918 +/- 0.006 (micro average: 0.918) (positive class: Positive) AUC (pessimistic): 0.839 +/- 0.012 (micro average: 0.839) (positive class: Positive) precision: 90.87% +/- 0.64% (micro average: 90.87%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3242 69 Positive: 611 6078 recall: 98.88% +/- 0.31% (micro average: 98.88%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3242 69 Positive: 611 6078 f_measure: 94.70% +/- 0.31% (micro average: 94.70%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3242 69 Positive: 611 6078		
SVM with SMOTE		
PerformanceVector: accuracy: 96.80% +/- 0.65% (micro average: 96.80%) ConfusionMatrix: True: Negative Positive Negative: 3793 260 Positive: 60 5887 AUC (optimistic): 0.994 +/- 0.002 (micro average: 0.994) (positive class: Positive) AUC: 0.994 +/- 0.002 (micro average: 0.994) (positive class: Positive) AUC (pessimistic): 0.994 +/- 0.002 (micro average: 0.994) (positive class: Positive) precision: 98.99% +/- 0.28% (micro average: 98.99%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3793 260 Positive: 60 5887 recall: 95.77% +/- 1.03% (micro average: 95.77%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3793 260 Positive: 60 5887 f_measure: 97.35% +/- 0.55% (micro average: 97.35%) (positive class: Positive) ConfusionMatrix: True: Negative Positive Negative: 3793 260 Positive: 60 5887		

Figure 8. Performance of NBC, DT, and SVM Using SMOTE

Figure 8 shows the NBC, DT, and SVM performance using SMOTE in Rapidminer. The performance results of the NBC model reveal notable metrics, with an accuracy of 86.98% \pm 0.96% and an AUC (Area Under the Curve) of 0.657 \pm 0.203, indicating its effectiveness in accurately classifying sentiments. The model demonstrates precision and recall rates of 100.00% and 78.82%, respectively, showcasing its ability to correctly identify positive sentiments while maintaining a low false positive rate. Additionally, the f-measure of 88.15% underscores the model's balanced performance in terms of precision and recall. These findings highlight the robustness and reliability of the NBC model in sentiment classification tasks, affirming its potential utility in analyzing and interpreting textual data for various applications. Otherwise, the performance evaluation of the Decision Tree (DT) model illustrates impressive metrics, with an accuracy of 93.20% \pm 0.42% and an AUC (Area Under the Curve) of 0.918 \pm 0.006, indicating its efficacy in sentiment classification tasks. Notably, the model demonstrates high precision and recall rates of 90.87% and 98.88%, respectively, underscoring its ability to identify positive sentiments while accurately minimizing false positives. Furthermore, the f-measure of 94.70% highlights the balanced performance of the DT model in terms of both precision and recall. These results affirm the robustness and reliability of the DT model in sentiment analysis, suggesting its suitability for various applications requiring accurate classification of textual data. In addition, the performance evaluation of the Support Vector Machine (SVM) model reveals impressive results, with an accuracy of 96.80% \pm 0.65% and an AUC (Area Under the Curve) of 0.994, indicating its effectiveness in sentiment classification tasks. Notably, the SVM model demonstrates high precision and recall rates of 98.99% and 95.77%, respectively, underscoring its ability to identify positive sentiments while minimizing false positives accurately. Furthermore, the f-measure of 97.35% highlights the balanced performance of the SVM model in terms of both precision and recall. These findings affirm the robustness and reliability of the SVM model in sentiment analysis, suggesting its suitability for various applications requiring accurate classification of textual data.

This study is confined to evaluating the performance of models or algorithms in sentiment data classification based on metrics such as accuracy, precision, recall, F-measure, and AUC. By focusing on these key performance indicators, the research aims to comprehensively assess the efficacy and reliability of the classification models in accurately predicting sentiment labels within the dataset. This rigorous evaluation approach enables a nuanced understanding of the model's strengths and limitations, facilitating informed decision-making regarding the suitability for sentiment analysis tasks. Consequently, adopting a systematic evaluation framework enhances the credibility and applicability of the study's findings in guiding future research and practical applications in sentiment analysis. The limitation of this research lies in both the dataset utilized and the framework employed, namely CRISP-DM. The dataset's constraints may impact the generalizability and comprehensiveness of the study's findings, potentially limiting the scope of insights derived from the analysis. Additionally, adopting the CRISP-DM framework while providing a structured approach to data mining tasks may impose constraints on the flexibility and adaptability of the research methodology to unique contextual factors. Furthermore, the research's scope neglects to address the livelihood and settlement contexts of the indigenous community of Kampung Naga, thereby overlooking crucial socio-economic and cultural dimensions that could enrich the analysis and broaden its applicability. Consequently, acknowledging and addressing these limitations is essential for ensuring the integrity and validity of the research findings and guiding future endeavors to address these gaps in knowledge.

4. CONCLUSION

Based on the findings derived from the analysis of sentiment classification models using the CRISP-DM methodology and data mining techniques, it is evident that the Naive Bayes Classifier (NBC), Decision Tree (DT), and Support Vector Machine (SVM) exhibit commendable performance in categorizing sentiments within digital content. The NBC model demonstrates an accuracy of 86.98% \pm 0.96%, precision of 100.00% \pm 0.00%, recall of 78.82% \pm 1.55%, and f-measure of 88.15% \pm 0.97%, with an AUC of 0.657 \pm 0.203. Likewise, the DT model showcases an accuracy of 93.20% \pm 0.42%, precision of 90.87% \pm 0.64%, recall of 98.88% \pm 0.31%, and f-measure of 94.70% \pm 0.31%, with an AUC of 0.918 \pm 0.006. Furthermore, the SVM model exhibits an accuracy of 96.80% \pm 0.65%, precision of 98.99% \pm 0.28%, recall of 95.77% \pm 1.03%, and f-measure of 97.35% \pm 0.55%, with an AUC of 0.994. These results underscore the robustness and reliability of these models in accurately classifying sentiments within textual data, thus affirming their potential utility in digital content analysis and data mining tasks. Moving forward, it is recommended to explore ensemble methods, continuously update the models with new data, investigate alternative sampling techniques, experiment with different feature engineering approaches, and collaborate with domain experts to enhance the effectiveness and applicability of sentiment analysis models in real-world scenarios.

ACKNOWLEDGEMENTS

Thanks to the Tourism Department, Faculty of Business Administration and Communication, and the Atma Jaya Catholic University of Indonesia.

REFERENCES

- [1] O. B. Onyancha, "Indigenous knowledge, traditional knowledge and local knowledge: what is the difference? An informetrics perspective," *Glob. Knowledge, Mem. Commun.*, vol. 73, no. 3, pp. 237–257, Jan. 2024, doi: 10.1108/GKMC-01-2022-0011.

- [2] Y. Yang, X. Lin, and R. B. Anderson, "Entrepreneurship by indigenous people in Canada and Australia: diverse modes and community implications," *Int. J. Entrep. Behav. Res.*, vol. 30, no. 1, pp. 90–109, Jan. 2024, doi: 10.1108/IJEBR-01-2023-0085.
- [3] B. Beamer and K. C. Gleason, "Reflections on the impact of informal sector tourism on indigenous Namibian Craft processes," *Arts Mark.*, vol. 12, no. 1, pp. 1–16, Jan. 2022, doi: 10.1108/AAM-05-2020-0015.
- [4] C. Makate, "Local institutions and indigenous knowledge in adoption and scaling of climate-smart agricultural innovations among sub-Saharan smallholder farmers," *Int. J. Clim. Chang. Strateg. Manag.*, vol. 12, no. 2, pp. 270–287, Jan. 2020, doi: 10.1108/IJCCSM-07-2018-0055.
- [5] L. Koppenhafer, K. Scott, T. Weaver, and M. Mulder, "The service empowerment model: a collaborative approach to reducing vulnerability," *J. Serv. Mark.*, vol. 37, no. 7, pp. 911–926, Jan. 2023, doi: 10.1108/JSM-10-2022-0317.
- [6] S. Lee, Y. Chang, O. K. D. Lee, S. Ryu, and Q. Yin, "Exploring online social platform affordances for digital creators: a multi-method approach using qualitative and configurational analysis," *Ind. Manag. Data Syst.*, vol. 124, no. 4, pp. 1501–1530, Jan. 2024, doi: 10.1108/IMDS-12-2023-0951.
- [7] N. Gryllakis and M. Matsiola, "Digital audiovisual content in marketing and distributing cultural products during the COVID-19 pandemic in Greece," *Arts Mark.*, vol. 13, no. 1, pp. 4–19, Jan. 2023, doi: 10.1108/AAM-09-2021-0053.
- [8] N. Nicoli, K. Henriksen, M. Komodromos, and D. Tsagalas, "Investigating digital storytelling for the creation of positively engaging digital content," *EuroMed J. Bus.*, vol. 17, no. 2, pp. 157–173, Jan. 2022, doi: 10.1108/EMJB-03-2021-0036.
- [9] A. P. Kieling, R. Tezza, and G. L. Vargas, "Website stage model for Brazilian wineries: an analysis of presence in digital and mobile media," *Int. J. Wine Bus. Res.*, vol. 35, no. 1, pp. 45–65, Jan. 2023, doi: 10.1108/IJWBR-05-2021-0032.
- [10] C. Clune and E. McDaid, "Content moderation on social media: constructing accountability in the digital space," *Accounting, Audit. Account. J.*, vol. 37, no. 1, pp. 257–279, Jan. 2024, doi: 10.1108/AAAJ-11-2022-6119.
- [11] C. Chen and T. Kellison, "The clock is ticking: contexts, tensions and opportunities for addressing environmental justice in sport management," *Sport. Bus. Manag. An Int. J.*, vol. 13, no. 3, pp. 376–396, Jan. 2023, doi: 10.1108/sbm-08-2022-0071.
- [12] R. Colbourne, P. Moroz, C. Hall, K. Lendsay, and R. B. Anderson, "Indigenous works and two eyed seeing: mapping the case for indigenous-led research," *Qual. Res. Organ. Manag. An Int. J.*, vol. 15, no. 1, pp. 68–86, Jan. 2020, doi: 10.1108/QROM-04-2019-1754.
- [13] A. Kaur and W. Qian, "The state of disclosures on Aboriginal engagement: an examination of Australian mining companies," *Meditari Account. Res.*, vol. 29, no. 2, pp. 345–370, Jan. 2020, doi: 10.1108/MEDAR-01-2020-0702.
- [14] L. Bellato and J. M. Cheer, "Inclusive and regenerative urban tourism: capacity development perspectives," *Int. J. Tour. Cities*, vol. 7, no. 4, pp. 943–961, Jan. 2021, doi: 10.1108/IJTC-08-2020-0167.
- [15] S. W. Maingi, "Safari tourism and its role in sustainable poverty eradication in East Africa: the case of Kenya," *Worldw. Hosp. Tour. Themes*, vol. 13, no. 1, pp. 81–94, Jan. 2021, doi: 10.1108/WHATT-08-2020-0084.
- [16] I. Moyo and H. M. S. Cele, "Protected areas and environmental conservation in KwaZulu-Natal, South Africa: on HEIs, livelihoods and sustainable development," *Int. J. Sustain. High. Educ.*, vol. 22, no. 7, pp. 1536–1551, Jan. 2021, doi: 10.1108/IJSHE-05-2020-0157.
- [17] P. Scherrer, "Tourism to serve culture: the evolution of an Aboriginal tourism business model in Australia," *Tour. Rev.*, vol. 75, no. 4, pp. 663–680, Jan. 2020, doi: 10.1108/TR-09-2019-0364.
- [18] T. Neha *et al.*, "Sustainable prosperity and enterprises for Maori communities in Aotearoa New Zealand: a review of the literature," *J. Enterprising Communities*, vol. 15, no. 4, pp. 608–625, Jan. 2021, doi: 10.1108/JEC-07-2020-0133.
- [19] R. Obiedat *et al.*, "Sentiment Analysis of Customers' Reviews Using a Hybrid Evolutionary SVM-Based Approach in an Imbalanced Data Distribution," *IEEE Access*, vol. 10, pp. 22260–22273, 2022, doi: 10.1109/ACCESS.2022.3149482.
- [20] Z. Wu, G. Cao, and W. Mo, "Multi-Tasking for Aspect-Based Sentiment Analysis via Constructing Auxiliary Self-Supervision ACOP Task," *IEEE Access*, vol. 11, no. May, pp. 82924–82932, 2023, doi: 10.1109/ACCESS.2023.3276320.
- [21] R. Bringula, S. A. I. D. A. Ulfä, J. P. P. Miranda, and F. A. L. Atienza, "Text mining analysis on students' expectations and anxieties towards data analytics course," *Cogent Eng.*, vol. 9, no. 1, 2022, doi: 10.1080/23311916.2022.2127469.
- [22] R. K. Botchway, A. B. Jibril, Z. K. Oplatková, and M. Chovancová, "Deductions from a Sub-Saharan African Bank's Tweets: A sentiment analysis approach," *Cogent Econ. Financ.*, vol. 8, no. 1, 2020, doi: 10.1080/23322039.2020.1776006.
- [23] R. Harakawa, T. Ogawa, and M. Haseyama, "Extracting Hierarchical Structure of Web Video Groups Based on Sentiment-Aware Signed Network Analysis," *IEEE Access*, vol. 5, pp. 16963–16973, 2017, doi: 10.1109/ACCESS.2017.2741098.
- [24] J. Chen, Q. Mao, and L. Xue, "Visual sentiment analysis with active learning," *IEEE Access*, vol. 8, pp. 185899–185908, 2020, doi: 10.1109/ACCESS.2020.3024948.
- [25] H. Zhang, S. Sun, Y. Hu, J. Liu, and Y. Guo, "Sentiment Classification for Chinese Text Based on Interactive Multitask Learning," *IEEE Access*, vol. 8, pp. 129626–129635, 2020, doi: 10.1109/ACCESS.2020.3007889.
- [26] C. B. Lee, H. N. Io, and H. Tang, "Sentiments and perceptions after a privacy breach incident," *Cogent Bus. Manag.*, vol. 9, no. 1, 2022, doi: 10.1080/23311975.2022.2050018.
- [27] V. Gupta, S. Singh, and S. S. Yadav, "The impact of media sentiments on IPO underpricing," *J. Asia Bus. Stud.*, vol. 16, no. 5, pp. 786–801, Jan. 2022, doi: 10.1108/JABS-10-2020-0404.
- [28] F. Cavigglioli, L. Lamberti, P. Landoni, and P. Meola, "Technology adoption news and corporate reputation: sentiment analysis about the introduction of Bitcoin," *J. Prod. Brand Manag.*, vol. 29, no. 7, pp. 877–897, Jan. 2020, doi: 10.1108/JPB-03-2018-1774.
- [29] W. Zheng, S. Zhang, C. Yang, and P. Hu, "Lightweight multilayer interactive attention network for aspect-based sentiment analysis," *Conn. Sci.*, vol. 35, no. 1, 2023, doi: 10.1080/09540091.2023.2189119.
- [30] K. Puh and M. Bagić Babac, "Predicting sentiment and rating of tourist reviews using machine learning," *J. Hosp. Tour. Insights*, vol. 6, no. 3, pp. 1188–1204, 2023, doi: 10.1108/JHTI-02-2022-0078.
- [31] F. Alattar and K. Shaalan, "A Survey on Opinion Reason Mining and Interpreting Sentiment Variations," *IEEE Access*, vol. 9, pp. 39636–39655, 2021, doi: 10.1109/ACCESS.2021.3063921.
- [32] J. Wu, K. Lu, S. Su, and S. Wang, "Chinese Micro-Blog Sentiment Analysis Based on Multiple Sentiment Dictionaries and Semantic Rule Sets," *IEEE Access*, vol. 7, pp. 183924–183939, 2019, doi: 10.1109/ACCESS.2019.2960655.
- [33] K. Xu, H. Zhao, and T. Liu, "Aspect-Specific Heterogeneous Graph Convolutional Network for Aspect-Based Sentiment Classification," *IEEE Access*, vol. 8, pp. 139346–139355, 2020, doi: 10.1109/ACCESS.2020.3012637.

- [34] I. Z. P. Hamdan and M. Othman, "Predicting Customer Loyalty Using Machine Learning for Hotel Industry," *J. Soft Comput. Data Min.*, vol. 3, no. 2, pp. 31–42, 2022.
- [35] Y. A. Singgalen, "Social Network Analysis and Sentiment Classification of Extended Reality Product Content," *J. Tek. Inform. C.I.T Medicom*, vol. 16, no. 1, pp. 24–34, 2024.
- [36] J. A. Syahid and D. Mahdiana, "Perbandingan algoritma untuk klasifikasi analisis sentimen terhadap Genose pada media sosial Twitter," *semantik*, vol. 7, no. 1, pp. 9–16, 2021, doi: 10.5281/zenodo.5034916.
- [37] H. Kim and G. Qin, "Summarizing Students' Free Responses for an Introductory Algebra-Based Physics Course Survey Using Cluster and Sentiment Analysis," *IEEE Access*, vol. 11, no. July, pp. 89052–89066, 2023, doi: 10.1109/ACCESS.2023.3305260.
- [38] K. Jahanbin and M. A. Z. Chahooki, "Aspect-Based Sentiment Analysis of Twitter Influencers to Predict the Trend of Cryptocurrencies Based on Hybrid Deep Transfer Learning Models," *IEEE Access*, vol. 11, no. November, pp. 121656–121670, 2023, doi: 10.1109/ACCESS.2023.3327060.
- [39] K. R. Mabokela, T. Celik, and M. Raborife, "Multilingual Sentiment Analysis for Under-Resourced Languages: A Systematic Review of the Landscape," *IEEE Access*, vol. 11, no. February, pp. 15996–16020, 2023, doi: 10.1109/ACCESS.2022.3224136.
- [40] T. Lin and I. Joe, "An Adaptive Masked Attention Mechanism to Act on the Local Text in a Global Context for Aspect-Based Sentiment Analysis," *IEEE Access*, vol. 11, no. May, pp. 43055–43066, 2023, doi: 10.1109/ACCESS.2023.3270927.