

Jeffery Dirden

ITAI 2373

July 26, 2025

NewsBot Reflection

Personal Collaboration Reflection

Since I completed this NewsBot project independently, all aspects of the work were my responsibility. While this meant I didn't have teammates to share ideas with or divide tasks, it also gave me the opportunity to develop a strong sense of ownership and self-discipline. I had to manage my time effectively, plan each step carefully, and troubleshoot problems without immediate external input. Though challenging at times, working solo helped me build confidence in my ability to handle complex projects from start to finish.

Dataset Selection and Technical Integration Challenges

Due to computing limitation, specifically running the project in Google Colab with limited memory and processing power. I carefully selected a dataset that balanced size and quality. I chose a subset of approximately 1,000 articles spread across four categories to ensure manageable processing while maintaining sufficient variety for classification. This dataset met the requirements for text richness and clear category labels, but its size was limited to prevent memory overload or long training times.

Integrating multiple technical components on my own was one of the most difficult parts of this project. Combining the TF-IDF features with sentiment scores and length metrics into a single feature set required careful data alignment. I encountered challenges with model training, especially tuning hyperparameters and managing data sparsity. Working with spaCy's NLP tools for entity extraction and dependency parsing also required multiple iterations to optimize performance and memory use. The absence of collaborators meant I had to rely heavily on research, testing, and self-debugging to overcome these hurdles.

Business Value Assessment

The NewsBot system I developed demonstrates significant potential business value. It automates the categorization and sentiment analysis of news articles, which can help media companies and marketers save time and make data-driven decisions.

Entity extraction and insights generation add depth, enabling quicker understanding of key people, organizations, and topics in the news. With solid classification accuracy, this system could support content recommendation and editorial monitoring, reducing manual effort and improving responsiveness to audience sentiment.

Individual Contributions

Since I worked alone, all contributions are mine. I designed and implemented the entire pipeline, including preprocessing, feature engineering, model training and evaluation, sentiment analysis, entity extraction, and insight generation. I integrated different data sources and methods to create a cohesive system based on the model we were provided. Additionally, I conducted comprehensive testing and visualization to ensure reliability. This solo effort strengthened my coding skills, project management, and problem-solving abilities throughout the entire machine learning workflow.

Future Enhancements

Looking forward, there are several ways I could enhance the system. Incorporating contextual embeddings like BERT would improve semantic understanding beyond TF-IDF. Adding advanced entity relationship and temporal analysis could offer richer insights. Addressing class imbalance with techniques like oversampling or ensemble learning could boost classification performance. Scaling the system for real-time processing and building an interactive dashboard would make it more practical for end users. These improvements would make NewsBot more powerful and user-friendly.

Professional Development Impact

This independent project has greatly contributed to my professional growth. It deepened my understanding of NLP and machine learning pipelines and gave me hands-on experience with real-world challenges in data integration and model tuning. Working solo strengthened my self-reliance, perseverance, and ability to learn new concepts independently. I enhanced my skills with tools like scikit-learn, spaCy, and data visualization libraries. Overall, this experience has prepared me well for future data science projects and increased my confidence in tackling complex technical problems independently.

Cited Sources

Bird, Steven, Edward Loper, and Ewan Klein. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. O'Reilly Media, 2009.

Honnibal, Matthew, and Ines Montani. *spaCy 101: Everything You Need to Know*. Explosion AI, <https://spacy.io>.

Pedregosa, Fabian, et al. "Scikit-learn: Machine Learning in Python." *Journal of Machine Learning Research*, vol. 12, 2011, pp. 2825–2830. <http://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html>.

Rehurek, Radim, and Petr Sojka. "Software Framework for Topic Modelling with Large Corpora." *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, ELRA, 2010. <http://is.muni.cz/publication/884893>.

"VADER Sentiment Analysis." *GitHub - Cjhutto/VaderSentiment*, <https://github.com/cjhutto/vaderSentiment>.

Waskom, Michael L. "Seaborn: Statistical Data Visualization." *Journal of Open Source Software*, vol. 6, no. 60, 2021, p. 3021. <https://doi.org/10.21105/joss.03021>.

"Matplotlib: Visualization with Python." *Matplotlib.org*, <https://matplotlib.org/>.

Google Colab. *Colaboratory*. Google Research, <https://colab.research.google.com/>.