

MLOps Tools for Deployment: A Case Study on Text Classification

Matea Lukić¹, Danica Ivković², Ana Poledica³

Abstract—This paper explores the integration of MLOps tools in the design and operationalization of machine learning pipelines for text classification. Focusing on a case study of web news classification, we examine the use of tools such as MLflow for experiment tracking, Docker for containerization, and Airflow for orchestration. The results reveal both promising advancements and significant limitations, underscoring the challenges of adopting DevOps practices in the rapidly evolving field of machine learning. Although the findings highlight the potential of MLOps to improve scalability and reproducibility, they also demonstrate that the domain is still in its early stages and requires further refinement.

Index Terms—MLOps, deployment, machine learning pipelines, text classification, scalability, reproducibility, operationalization

I. INTRODUCTION

In recent years, machine learning (ML) has evolved into a critical component of modern applications [2], from the recommendation systems to the ability to process natural language. However, deploying ML models reliably at scale remains a significant challenge [4], particularly because much of machine learning today still relies on academic datasets and tools that lack scalability for production [1].

MLOps, inspired by DevOps principles, addresses these issues by introducing automation, agility, and collaboration into the development of ML systems [3]. This discipline has emerged as a key to managing the complexities of the development and deployment of machine learning solutions [5]. Unlike traditional software engineering, ML systems pose unique challenges due to their reliance on data, iterative experimentation, and model lifecycle management [6].

The remainder of the paper is organized as follows. Section II summarizes the findings of the literature review. Section III outlines the development steps and discusses various tool options. Section IV presents a case study on an NLP project focusing on the construction of a multi-step ML pipeline for news classification. Finally, Section V concludes the paper with a summary and discusses future directions.

II. RELATED WORK

A. Overview of MLOps

Recent research sheds light on the evolution, components, and practices within the MLOps domain. Platforms like Mod-

elOps, TensorFlow Extended (TFX), and Kubeflow provide end-to-end lifecycle management, orchestrating stages such as data preparation, model training, validation, and deployment into comprehensive ML pipelines [2], [6]. However, challenges remain in resource optimization and efficient handling of pipeline stages. Key points that need to be addressed are bottlenecks such as GPU utilization inefficiencies, which could impact overall system performance [6], and issues like data distribution shifts and model retraining failures [2].

B. Problems in ML System Development

As explained in [7] ML systems often accumulate hidden technical debt, making them fragile and difficult to maintain. A significant issue is entanglement, where tightly coupled components cause small changes to propagate unpredictably, reflecting the CACE principle (Changing Anything Changes Everything). Common antipatterns make these issues even worse, such as glue code, where ad hoc scripts connect disparate system components, and pipeline jungles, complex, undocumented workflows that are hard to debug or reproduce.

C. Case Studies Based on MLOps Paradigm

The study [8] illustrates the effective use of MLOps principles in time-series forecasting. It employs tools such as Docker for containerization, Jenkins for CI/CD, unit tests for validation, PyTorch for model development, MLflow for experiment tracking, and PostgreSQL for data management.

Expanding these ideas [9] integrated tools like Dask, Katib, PyTorch Operator, and KServe within a single Kubeflow Pipelines (KFP) workflow. This setup allowed for anomaly detection using telemetry data from the International Space Station (ISS), showcasing how MLOps can handle real-world challenges by integrating various tools.

These studies provide a blueprint for implementing machine learning solutions in complex systems, demonstrating how MLOps principles and a variety of tools can be integrated to address real-world challenges while ensuring reproducibility, scalability, and automation.

III. ML DEVELOPMENT STEPS

This paper discusses the deployment of an ML project using the Python ecosystem. It highlights the need for tools to ensure that the solutions are functional in production. Python's versatility and libraries make it ideal for this use case. Successful deployment requires a clear, step-by-step process through all project stages.

¹Matea Lukić, Faculty of Organizational Sciences, University of Belgrade, Jove Ilića 154, 11000 Belgrade, Serbia (e-mail: lukimatea166@gmail.com)

²Danica Ivković, Factory World Wide, Bulevar Mihajla Pupina 115a, 11000 Belgrade, Serbia (e-mail: danica.ivkovic@factoryww.com)

³Ana Poledica, Faculty of Organizational Sciences, University of Belgrade, Jove Ilića 154, 11000 Belgrade, Serbia (e-mail: ana.poledica@fon.bg.ac.rs)

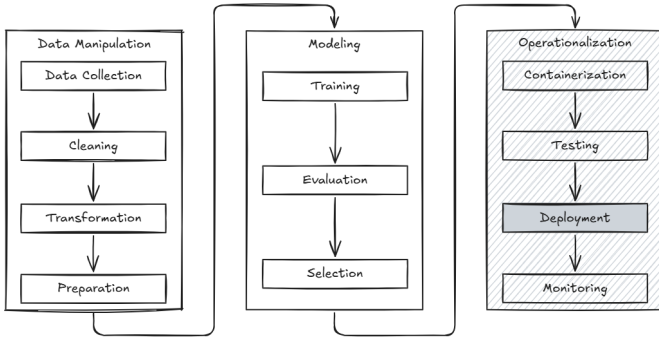


Fig. 1. Development steps in ML projects

As highlighted in [10], finding the right balance between simplicity and flexibility is crucial due to the wide range of available tools, from general-purpose to specialized solutions. The challenge lies in minimizing the complexity of the tool while ensuring adequate functionality. To explore this principle, the development steps [10], [14], are outlined in the Fig. 1 and will be explained further.

A. Data Manipulation pipeline

Data manipulation begins with acquiring the necessary data, which ideally already exist in structured repositories, but more often require collection and labeling. There are various methods for gathering data, from web scraping and using APIs to manually collecting it [3]. Platforms like Kaggle, UCI Machine Learning Repository and Hugging Face Datasets serve as valuable recourses for obtaining curated datasets, while libraries like Scrapy and BeautifulSoup simplify the process of web scraping. For cases where labeled data is needed, tools such as Label Studio, Prodigy, and Doccano assist in efficient annotation.

Data cleaning involves resolving issues such as duplicates, outliers, missing values, and inconsistencies, whereas data transformation involves altering the type or distribution of variables, such as converting numeric values into categorical ones, normalizing the data or even creating new derived variables to enhance model performance [11]. Libraries like Pandas and Polars are widely used for data cleaning and transformation, providing extensive tools for data manipulation. Missing values can be imputed using techniques like mean substitution or K-Nearest Neighbors (KNN) with libraries such as Scikit-learn, while PyOD offers methods for identifying and managing outliers.

Data preparation, as outlined in [14], involves organizing the dataset for training and evaluation. This process typically includes splitting the data into training, validation, and test subsets using methods such as random splitting or stratified sampling, with libraries like Scikit-learn providing built-in support. In addition to splitting, the concept of spanning allows for the inclusion of new data snapshots over time, ensuring that models remain up-to-date without having to retrain from scratch. To maintain reproducibility and consistency, tools like Data Version Control or Pachyderm can be used.

B. Modeling pipeline

Training is the first step in modeling pipeline, where a model learns patterns from data by adjusting its parameters [12]. Modern frameworks like TensorFlow, JAX and PyTorch simplify and enhance this process, with PyTorch Lightning providing a streamlined interface for PyTorch. Platforms like Google Colab are essential for providing cloud-based training environments with free GPU/TPU access, making it possible for researchers and practitioners to train computationally intensive models without investing in expensive hardware. Jupyter Notebook and Kaggle Kernels complement this by offering options for both local and cloud-based experimentation, allowing for rapid prototyping, visualization, and interactive debugging. Tools like MLflow, Weights & Biases (Wandb), and Neptune.ai track experiments and ensure reproducibility by recording hyperparameters, metrics, and artifacts.

Evaluation is critical to understanding a model's performance and ensuring its generalizability. Tools like scikit-learn offer a wide range of metrics for evaluating models in classification, regression, and clustering tasks. SciPy and StatsModels provide additional statistical tools for deeper analysis and hypothesis testing. For visualizing and comparing model performance, platforms such as DAGsHub, Guild AI, and Comet.ml are widely used. These tools facilitate identifying overfitting, underfitting, and other performance bottlenecks. The ONNX (Open Neural Network Exchange) format allows evaluation and benchmarking across different frameworks, ensuring interoperability. Additionally, models can be converted to other formats such as H5, ProtoBuf, Pickle, Joblib and TFJS enabling compatibility with specific deployment environments.

Selection involves choosing the best architecture and hyperparameters to achieve the optimal level of flexibility for a model [12]. Tools like Optuna, Ray Tune, and Hyperopt provide automated hyperparameter optimization, leveraging techniques like Bayesian optimization and distributed execution. Experiment management platforms such as MLflow enable tracking and comparing multiple model versions, aiding informed decision-making. For NLP tasks, libraries like Hugging Face (HF) Transformers, spaCy, and AllenNLP offer pre-trained models that simplify selection and customization. Cloud-based platforms like AWS SageMaker, Google AI Platform, and Azure Machine Learning provide managed services for large-scale model selection and deployment.

C. Operationalization pipeline

Containerization is a key concept in modern software development, where applications are packaged in isolated environments called containers. Containers are standalone and executable packages of software that include everything needed to run an application, such as code, system tools, libraries, and settings [13]. While Docker is one of the most widely used tools for this process, Podman offers a daemonless alternative with enhanced security features, such as rootless containers. Kubernetes (k8s), on the other hand, is not a container engine but a container orchestration platform

designed to manage and automate the deployment, scaling, and operation of containerized applications across clusters. Unlike Docker and Podman, which focus on creating and running individual containers, Kubernetes operates at a higher level, coordinating the behavior of multiple containers and ensuring high availability and scalability.

The next step is testing, which ensures the reliability and robustness of machine learning pipelines. Libraries such as pytest and unittest support unit testing of pipeline components, while Great Expectations and Deequ validate the quality of input and output data. For integration and system testing, frameworks such as Airflow, Prefect, and Kubeflow Pipelines manage end-to-end workflows, ensuring seamless execution across all stages. Tools like Jenkins provide automated continuous integration and testing pipelines.

Deployment transforms machine learning models into production ready solutions [1]. Tools like AWS SageMaker, Google AI Platform, and Azure Machine Learning provide managed services for scalable deployment. Open-source frameworks like FastAPI, Flask, and Django facilitate the creation of custom APIs for model inference. Exporting models in formats such as ONNX and TFJS ensures compatibility with edge devices and mobile applications.

Monitoring ensures that deployed models continue to perform as expected. Tools like Prometheus and Grafana provide real-time metrics and visualization, enabling proactive issue resolution. Platforms like Evidently AI, WhyLabs, and Fiddler monitor data drift (the delta between the changes in the data from the last time the model training occurred [1]), model performance, and fairness, ensuring models remain reliable over time. Cloud-native monitoring solutions, such as Amazon CloudWatch and Google Cloud Monitoring, integrate seamlessly with managed deployment services, providing holistic oversight. Logging tools like ELK Stack (Elasticsearch, Logstash, Kibana) and Fluentd capture detailed logs for troubleshooting and auditing.

IV. CASE STUDY: TEXT CLASSIFICATION

This case study demonstrates the development and operationalization of a machine learning system designed for daily classification of scraped news using the AG News dataset¹. The workflow encompassed data collection, preprocessing, model training, evaluation, and deployment, leveraging a variety of modern tools and frameworks, shown in Fig. 2, to ensure robustness and efficiency. Table I summarizes the benefits and alternatives of these tools, while some of their limitations will be discussed later in the text.



Fig. 2. Tools Used in Each Development Step

TABLE I
OVERVIEW OF TOOLS, BENEFITS, AND ALTERNATIVES

Tool	Benefits	Alternatives
BeautifulSoup	Fast, lightweight	Scrappy, Selenium
Polars	Multithreaded, scalable	Pandas, Spark, Dask
Scikit-learn	Documentation, popular	Orange, RapidMiner
PyTorch	Syntax, modern, intuitive	TensorFlow, JAX
Tensorflow	TensorBoard, scalable	PyTorch, JAX
MLFlow	Reproducibility, flexibility	Wandb, Neptune.ai
Google Colab	Fast, Google Drive access	Kaggle, Local build
HF Transformers	SOTA models	AllenNLP, SpaCy
DagsHub	Github integration	Databricks, Wandb
ONNX	Standard model format	H5, TFJS, Joblib
Matplotlib	Most popular	Seaborn, Plotly
Airflow	Wide set of providers	Dagster, Prefect
Docker	Ecosystem, easy to learn	Podman, K8s
FastApi	Less code, performance	Flask, Django
Uvicorn	ASGI server, lightweight	Daphne, Hypercorn
PostgreSQL	Scalable, cheap	SQLite, Casandra

A. Methodology Used

The initial steps of the pipeline involved data collection and preprocessing. The AG News dataset, sourced from the Kaggle platform², was used for training and testing. Once trained, the model was applied to classify news articles, which were scraped daily from various online sources and parsed using BeautifulSoup. Project package management was handled using Poetry.

The data manipulation process involved multiple preprocessing steps to standardize textual data for model training. This included removing HTML tags, converting the text to lowercase for consistency, cleaning the text by eliminating references to news agency names, and combining the title and the news description into a single, coherent text field. Polars was used to efficiently handle these tasks.

The modeling phase involved training a variety of models using different frameworks and libraries. PyTorch was used to implement Recurrent Neural Networks (RNNs) [15] and Long Short-Term Memory (LSTM) networks [16]. For transformer-based models, Hugging Face’s ELECTRA [17] and DistilBERT [18] were utilized, leveraging pre-trained architectures. The text was tokenized with TensorFlow, which, along with other preprocessing steps, prepared the data for model input. Scikit-learn was used to compute the evaluation metrics necessary for comparing the performance of different models. These metrics, including accuracy, precision, recall, and F1-score, were then visualized using Matplotlib to provide clear insights into model performance and guide further improvements.

The combined use of these frameworks facilitated a comprehensive exploration of modeling techniques. Model training was performed on Google Colab, taking advantage of its free GPU and TPU resources to accelerate computations. Experiment tracking, including hyperparameter configurations and performance metrics, was managed through MLflow, which was integrated with DAGsHub to ensure reproducibility and

¹Project Repository: <https://github.com/MateaLukiccc/MLOps-For-NLP>

²<https://www.kaggle.com/datasets/amananandrai/ag-news-classification-dataset>

efficient collaboration. The final model was exported in ONNX format, ensuring cross-platform compatibility and optimized inference.

The operationalization of the system was designed to ensure reliability, scalability, and automation. FastAPI exposed the trained model for real-time inference, providing a fast and efficient API endpoint, while Uvicorn served as the asynchronous web server, enabling quick handling of classification requests. Docker was used for containerization, packaging the application along with its dependencies. Task scheduling was orchestrated with Airflow, automating daily scraping and classification processes, which were set to run at midnight. The classification results were stored in PostgreSQL, providing efficient querying and retrieval capabilities for downstream applications.

B. Observations

Integrating traditional DevOps tools into machine learning workflows presented several challenges. Poetry falls short in addressing specific needs in machine learning workflows, particularly when it comes to specifying the 'CPU-only' version of the Torch library. This limitation complicates workflows when models are trained on platforms like Google Colab rather than on machines with local GPUs.

Apache Airflow, on the other hand, introduced a different set of challenges when deployed in a Docker Compose environment. The official documentation³ explicitly states that this setup is suitable for learning and experimentation but not recommended for production systems due to the lack of security guarantees. Adapting it for real-world use requires specialized expertise in Docker and Docker Compose, making it less accessible and limiting the support available from the Airflow community.

These examples highlight that MLOps practices are still evolving to meet the demands of machine learning workflows. Although these tools are widely used and well-established in the DevOps world, their integration with machine learning projects requires additional adjustments and workarounds.

V. CONCLUSION

This paper provides a comprehensive overview of MLOps tools and their practical application in the context of text classification, a common NLP problem. The case study on news classification using the AG News dataset highlighted the importance of integrating and automating workflows to enhance model development and deployment. Tools like PyTorch and Hugging Face Transformers were employed for experimenting with different model architectures, while the primary focus was on leveraging MLOps tools such as MLflow for tracking model performance, and Docker and Airflow for containerizing the models and automating the classification pipeline.

The practical application of MLOps in NLP showed their potential to increase the efficiency, reproducibility, and scalability of machine learning workflows. These frameworks

enable practitioners to address common challenges, such as managing numerous experiments, handling data distribution shifts, and ensuring real-time monitoring. However, integrating these components into a cohesive MLOps pipeline remains a significant challenge, underscoring that the MLOps field is still in its early stages. As the field evolves, staying informed about emerging tools and best practices will be crucial for effectively managing the complexities of NLP problems and driving continued innovation.

REFERENCES

- [1] Gift, N., & Deza, A. (2021). Practical MLOps. " O'Reilly Media, Inc."
- [2] Wazir, S., Kashyap, G. S., & Saxena, P. (2023). Mlops: A review. arXiv preprint arXiv:2308.10908.
- [3] Haviv, Y., & Gift, N. (2023). Implementing MLOps in the Enterprise. " O'Reilly Media, Inc."
- [4] Kreuzberger, D., Kühl, N., & Hirschl, S. (2023). Machine learning operations (mlops): Overview, definition, and architecture. IEEE access, 11, 31866-31879.
- [5] Alla, S., Adari, S. K., Alla, S., & Adari, S. K. (2021). What is mlops?. Beginning MLOps with MLFlow: Deploy Models in AWS SageMaker, Google Cloud, and Microsoft Azure, 79-124.
- [6] Zhou, Y., Yu, Y., & Ding, B. (2020, October). Towards mlops: A case study of ml pipeline platform. In 2020 International conference on artificial intelligence and computer engineering (ICAICE) (pp. 494-500). IEEE.
- [7] Sculley, D., Holt, G., Golovin, D., Davydov, E., Phillips, T., Ebner, D., ... & Dennison, D. (2015). Hidden technical debt in machine learning systems. Advances in neural information processing systems, 28.
- [8] Subramanya, R., Sierla, S., & Vyatkin, V. (2022). From DevOps to MLOps: Overview and application to electricity market forecasting. Applied Sciences, 12(19), 9851.
- [9] Steude, H. S., Geier, C., Moddemann, L., Creutzenberg, M., Pfeifer, J., Turk, S., & Niggemann, O. (2024). End-to-end MLOps integration: a case study with ISS telemetry data. In ML4CPS—Machine Learning for Cyber-Physical Systems.
- [10] Symeonidis, G., Nerantzis, E., Kazakis, A., & Papakostas, G. A. (2022, January). Mlops-definitions, tools and challenges. In 2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC) (pp. 0453-0460). IEEE.
- [11] Brownlee, J. (2020). Data preparation for machine learning: data cleaning, feature selection, and data transforms in Python. Machine Learning Mastery.
- [12] James, G. (2013). An introduction to statistical learning.
- [13] Merkel, D. (2014). Docker: lightweight linux containers for consistent development and deployment. Linux j, 239(2), 2.
- [14] Hapke, H., & Nelson, C. (2020). Building machine learning pipelines. O'Reilly Media.
- [15] Schmidt, R. M. (2019). Recurrent neural networks (rnns): A gentle introduction and overview. arXiv preprint arXiv:1912.05911.
- [16] Hochreiter, S. (1997). Long Short-term Memory. Neural Computation MIT-Press.
- [17] Clark, K. (2020). Electra: Pre-training text encoders as discriminators rather than generators. arXiv preprint arXiv:2003.10555.
- [18] Sanh, V. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108.

³Airflow Documentation: <https://airflow.apache.org/>