

Project Final Report

Group 15

Sanjana Shailesh Bongale (675620038)

Divya Pathak (667449971)

Table of Contents

1. Introduction.....	2
1.1. Project Overview	2
1.2. Rationale for the Project and Tools Explored	2
1.3. Exploratory Data Analysis (EDA).....	2
1.4. Experiment Tracking.....	2
1.5. Role of various MLOps Tools in Streamlining Machine Learning Workflows	2
2. Overview and Features of the Tools Used.....	3
2.1. DataPrep.....	3
Features of DataPrep.....	3
2.2. Neptune.ai.....	3
Features of Neptune.ai	3
3. Alternative Tools.....	4
3.1. Comparison of a few EDA Tools with DataPrep.....	4
3.2. Comparison of a few Experiment Tracking Tools with Neptune.ai	5
4. Key Steps in Project Methodology.....	6
Step 1 - Setting up Neptune.ai for the team	6
Step 2 - Installation of Required Libraries	6
Step 3 - Performing EDA using DataPrep	6
Step 4 - Performing Experiment Tracking using Neptune.ai	7
Step 5 - Uploading Visualizations to Neptune.ai.....	7
Step 6 - Stopping the Run	7
Step 7 - On Neptune.ai	7
5. Lessons Learned.....	7
5.1. DataPrep.....	7
5.2. Neptune.ai.....	8
5.3. Key Learning Points from the Project.....	8
6. Conclusion.....	9
7. Project Attachments	9

1. Introduction

1.1. Project Overview

This project focuses on extracting valuable insights from any dataset through exploratory data analysis (EDA) with the help of the DataPrep library. Experiment tracking is facilitated using Neptune.ai. For this study, we applied these techniques to the Titanic dataset.

1.2. Rationale for the Project and Tools Explored

This project was undertaken as a group effort to develop our skills in data analysis and gain practical experience with MLOps tools. Our primary focus was to explore the DataPrep library, which allowed us to perform exploratory data analysis (EDA) efficiently. We utilized its key functions to load the dataset, generate a comprehensive EDA report, and perform visualization. Additionally, we integrated Neptune.ai for experiment tracking, enabling us to log the changes made during each version systematically. This hands-on experience with these tools provided us with foundational knowledge in data preparation and experiment management as beginners in the field.

1.3. Exploratory Data Analysis (EDA)

The process of data analysis consists of an essential phase of exploratory data analysis (EDA), which is a technique that is used to summarize and visualize datasets to detect patterns, trends, and anomalies. It does so by looking at the data distributions and exploring the relationships between variables with the help of a range of statistical techniques and graphical tools, such as histograms, scatter plots, and descriptive statistics.

EDA directs additional analysis by assisting in the discovery of possible data problems, such as missing values and outliers. Moreover, EDA enables data scientists to create hypotheses and offer insights that help in the selection of the most effective method for the data-driven decision-making process by enabling a deeper comprehension of the data.

1.4. Experiment Tracking

Experiment tracking is a systematic approach for monitoring machine learning experiments, encompassing parameters, measurements, and results. This practice ensures transparency and reproducibility in model development. Experiment tracking tools enable data scientists to log and visualize critical information, such as model configurations, training data versions, and evaluation metrics.

Version control plays a pivotal role in this process, allowing teams to manage different datasets and model iterations. It ensures that every component can be traced back to its origin, thereby enhancing collaborative efforts. By streamlining the comparison of outcomes and identifying the best-performing models, this systematic approach fosters a more efficient and reliable machine-learning workflow.

1.5. Role of various MLOps Tools in Streamlining Machine Learning Workflows

- **Data Preparation Tools:** Automate the data cleaning and exploratory data analysis process, ensuring data readiness with minimal manual effort.
- **Model Development Tools:** Facilitate version control for code changes, improving collaboration and organization of model iterations.
- **Model Training Tools:** Monitor training runs and manage hyperparameters, allowing for real-time visualization of performance metrics.
- **Model Deployment Tools:** Automate CI/CD pipelines for efficient integration of models into production environments.
- **Monitoring and Maintenance Tools:** Track the model performance in real time to detect drift or anomalies for timely adjustments.
- **Collaboration and Reporting Tools:** Enhance communication through shared dashboards and reports, streamlining team collaboration.

2. Overview and Features of the Tools Used

2.1. DataPrep

DataPrep is a Python library specifically designed to simplify data preparation and exploratory data analysis (EDA). Aimed at data scientists and analysts, DataPrep streamlines the process of preparing data for analysis by providing intuitive functions that automate common tasks.

The library allows users to easily load datasets, perform data cleaning, and generate insightful visualizations with minimal coding effort. With its user-friendly interface, DataPrep enables beginners to conduct comprehensive EDA without extensive programming knowledge.

Features of DataPrep

DataPrep offers several key features that significantly enhance the data preparation and exploratory data analysis (EDA) process:

- **Automated EDA Report Generation:** DataPrep allows users to create comprehensive Exploratory Data Analysis (EDA) reports with a single function call, `create_report()`. This function provides a detailed summary of the dataset, including key insights for numeric, categorical, and datetime variables. The report visualizes critical statistics such as missing values, data distributions, correlations, and outliers. It helps users quickly understand the underlying structure and relationships within the data by automatically generating summary reports for each type of variable.
- **Visualization Functions:** DataPrep provides a range of built-in visualization tools, making it easy to create charts and plots. With minimal code, users can generate histograms, scatter plots, box plots, bar charts, and correlation heatmaps to explore relationships and data distributions using the `plot()` function. These ready-to-use visualizations simplify the process of gaining insights from the data.
- **Data Cleaning and Preparation:** The `DataPrep.clean()` function efficiently handles missing values, duplicates, and improper data formats and offers specialized cleaning tools for text, emails, phone numbers, and IP addresses. This makes it easier to prepare datasets for analysis.
- **User-Friendly Interface:** Designed for both beginners and experienced users, DataPrep's intuitive functions allow users to focus on insights rather than complex coding, enhancing productivity in the data analysis workflow.

These features collectively make DataPrep a valuable tool for streamlining the EDA process and improving data-driven decision-making.

2.2. Neptune.ai

Neptune.ai is a powerful platform designed for experiment tracking and model monitoring in machine learning workflows. It provides data scientists and machine learning engineers with a centralized solution to log, organize, and visualize their experiments, making it easier to track performance metrics and compare different model iterations.

Neptune.ai enhances collaboration within teams by allowing users to share insights and results seamlessly. The platform supports version control for datasets and models, ensuring that all team members work with the correct components and can trace back to previous versions. With its intuitive user interface, Neptune.ai empowers users to focus on developing high-quality models while maintaining a clear record of their experimentation process, ultimately streamlining the path from model development to deployment.

Features of Neptune.ai

Neptune.ai offers several key features that significantly enhance the management of machine learning experiments:

- **Experiment Tracking:** Neptune.ai allows users to log parameters, hyperparameters, and performance metrics across multiple experiments. This helps in keeping track of all experiments, making comparisons, and debugging easy.

- **Visualization & Dashboard:** Provides real-time visualization of metrics, including loss and accuracy. Users can create custom dashboards to monitor performance trends and experiment results, enhancing interpretability and insights.
- **Model Management:** Neptune.ai supports model versioning, enabling users to save, compare, and manage multiple model versions. This ensures easy access to past models and seamless experimentation with different iterations.
- **Collaboration:** Facilitates collaboration by allowing teams to share experiments, results, and visualizations. Multiple contributors can access the same projects, enhancing teamwork and collective decision-making.
- **Artifact Tracking:** Automatically tracks and stores artifacts like datasets, model weights, and visualizations, ensuring comprehensive documentation. This is crucial for experiment reproducibility and versioning.
- **API & SDK Integration:** Neptune.ai provides Python and REST APIs for seamless integration into workflows. These allow custom logging, automated tracking, and flexible experiment management for more sophisticated projects.

These features collectively make Neptune.ai an invaluable tool for managing the complexities of machine learning workflows and fostering a culture of experimentation and collaboration within teams.

3. Alternative Tools

In the rapidly evolving landscape of data science, various tools are available for Exploratory Data Analysis (EDA) and experiment tracking. This section outlines some notable alternatives to DataPrep and Neptune.ai, highlighting their features and unique advantages. We also provide comparisons and reasons for choosing DataPrep and Neptune.ai over these alternatives based on our experiences.

3.1. Comparison of a few EDA Tools with DataPrep

Feature/Aspect	DataPrep	Apache Zeppelin	Pandas Profiling
Primary Functionality	Simplifies data preparation and EDA with automated report generation.	Interactive data analytics and collaborative documents for multiple languages.	Generates HTML reports for quick insights from pandas DataFrames.
User Experience	User-friendly, designed for Python beginners; focuses on simplicity.	More complex interface, suitable for collaborative and multi-language analytics.	Simple to use for generating reports but limited to pandas DataFrames.
Integration	Integrates seamlessly with popular Python libraries.	Supports various data processing languages like Python, Scala, and R; requires more setup.	Directly integrates with pandas, making it easy for Python users.
Visualizations	Provides automated visualizations as part of the report.	Offers a range of visualization tools but may require more manual setup.	Generates visualizations as part of the HTML report but is less interactive.
Collaboration	Limited collaboration features focused on individual use.	Built for collaboration with shared notebooks and live editing capabilities.	Primarily focused on generating reports, not designed for collaboration.

DataPrep was chosen for this project due to its simplicity and ease of use, making it an ideal tool for beginners in data analysis. As our group aimed to familiarize ourselves with exploratory data analysis (EDA) processes, DataPrep's intuitive functions allowed us to quickly load and visualize the Titanic dataset and generate a report.

The library’s capability to generate automated EDA reports streamlined our workflow, enabling us to focus on interpreting the data. This user-friendly approach not only facilitated our learning experience but also helped us effectively communicate insights and study the visualizations derived from the data, enhancing our overall understanding of the analysis process.

3.2. Comparison of a few Experiment Tracking Tools with Neptune.ai

Feature/Aspect	Neptune.ai	Weights and Biases (W&B)	Comet
Experiment Tracking	User-friendly interface to log and visualize metrics easily.	Rich visualizations and detailed metrics tracking; strong collaboration features.	Comprehensive dashboard for monitoring experiments and dataset tracking.
User Experience	Lightweight and intuitive, suitable for beginners.	User-friendly but may have a steeper learning curve for new users.	Feature-rich but may feel cluttered for beginners.
Integration	Integrates well with Google Colab and various tools for a smooth workflow.	Broad compatibility with numerous frameworks but may require more configuration.	Extensive dataset management and integration options.
Data Management	Focuses primarily on experiment tracking with some visualization capabilities.	Offers extensive data tracking along with experiment management.	Strong emphasis on tracking datasets, code changes, and experiments.
User Interface	Clean, minimalistic design that helps users focus on experiments.	Rich interface but may appear cluttered due to many features.	Feature-rich interface that may overwhelm beginners.

Neptune.ai was chosen for our project due to its user-friendly interface and comprehensive features tailored for experiment tracking and visualizations, making it particularly suitable. While other alternatives offer valuable functionalities, Neptune.ai stands out in several ways.

Overall, Neptune.ai provided the ideal balance of functionality and simplicity for our group, allowing us to focus on developing our skills in experiment tracking and data visualization.

4. Key Steps in Project Methodology

In our project, we followed a systematic approach to set up the environment, install required libraries, and utilize tools for data analysis and experiment tracking. Google Colab was used in the project to write and execute our code efficiently.

Step 1 - Setting up Neptune.ai for the team

- 1: Creating a Neptune.ai Account:** We started by visiting neptune.ai and signed up using our email IDs.
- 2: Setting Up a Workspace:** Next, we created a workspace by selecting "Create a Workspace" and naming it. To collaborate effectively, we added both team members by going to the "Workspace Settings" and using the "Invite Member" option.
- 3: Creating a New Project:** Within the workspace, we created our project by clicking "New Project", providing a descriptive name, and selecting the relevant key and privacy settings. The privacy settings are kept to Public for this project.
- 4: Generating an API Token:** We generated an API token by navigating to Profile > Get your API, which was used to authenticate Neptune from our code.

Step 2 - Installation of Required Libraries

The necessary libraries were installed to facilitate experiment tracking, data exploration, and visualizations.

Command:

```
!pip install neptune-client dataprep seaborn matplotlib
```

Purpose of Libraries:

Neptune.ai: Used for tracking experiments and visualizing results.

DataPrep: Simplifies the process of exploratory data analysis (EDA) through automated reporting.

Seaborn & Matplotlib: Libraries for creating customized visualizations.

Step 3 - Performing EDA using DataPrep

After installing the required libraries, we utilized DataPrep to load the Titanic dataset and generate an EDA report. This process included:

- **Loading the Titanic Dataset:** We used the `load_dataset()` function from the DataPrep library to easily import the Titanic dataset. This straightforward method saved us time and effort compared to manually loading and preparing the data.
- **Generating an EDA Report:** By calling the `create_report(df)` function, we quickly generated an in-depth EDA report. This report included key statistics, visualizations, and insights about the dataset, such as missing values and data distributions. It provided an efficient way to understand our data at a glance.
- **Visualizing Key Variables:** The `plot()` function was used to visualize the Age variable in the dataset. This allowed us to see how age was distributed among the passengers, providing insights into the demographic composition of the Titanic's travelers.

Sample Code:

To Create Report->

```
df = load_dataset("titanic")  
create_report(df)
```

To Save Report->

```
report = create_report(df)  
report.save("titanic_eda_report.html")
```

To Download Report->

```
from google.colab import files  
files.download("titanic_eda_report.html")
```

Step 4 - Performing Experiment Tracking using Neptune.ai

Neptune.ai is integrated into the project to effectively track and manage our EDA experiments and visualizations. This process included. To initialize a run in Neptune.ai, the `neptune.init_run()` function is used. The project details and API token generated are mentioned in this function. This step establishes a connection between the code and the Neptune platform.

Sample Code:

```
run = neptune.init_run(  
    project='workspace_name/project_name',  
    api_token='api_token_generated_by_neptune.ai_in_step_0'  
)
```

Step 5 - Uploading Visualizations to Neptune.ai

In this step, the custom visualizations created using Matplotlib and Seaborn were uploaded to Neptune.ai. This made the visual outputs easily accessible for future reference.

Sample Code:

```
run['visualizations/titanic_age_distribution_chart'].upload('visualization_image_name.png')
```

Step 6 - Stopping the Run

The final step is to stop the run with the `run.stop()` function. This step finalizes and saves all the logged information and visualizations to the mentioned Neptune project.

Step 7 - On Neptune.ai

On Neptune.ai, the below observations and analysis were made by our team:

- A detailed list of all experiment runs, along with timestamps, was displayed for tracking and project organization.
- Side-by-side comparisons of experiment runs were available, highlighting changes in parameters and metrics for easy analysis.
- A visual display compared images generated from each run, enabling quick evaluation of results across different experiments.
- A custom visualization comparing female vs male data across 2 different runs was created on the dashboard, allowing gender-based insights through multiple experiment analyses.

5. Lessons Learned

5.1. DataPrep

Using DataPrep significantly enhanced our understanding of automated Exploratory Data Analysis (EDA). The library's ability to quickly generate comprehensive reports and visualizations allowed us to grasp insights without getting bogged down in the nitty-gritty of data manipulation. This experience reinforced the idea that automating certain EDA tasks can save time and effort, especially for beginners who may struggle with the intricacies of data analysis. Below are some advantages and disadvantages of using DataPrep for EDA:

Advantages	Disadvantages
Ease of generating insightful reports with minimal code, making it accessible for beginners in data science.	Lacks flexibility compared to manual approaches using libraries like Pandas and Seaborn.
Saves time by automating the reporting process, enhancing productivity.	Standard visualizations may not capture specific messages or insights needed for more detailed analysis.

Advantages	Disadvantages
User-friendly interface that simplifies the exploratory data analysis process.	Manual crafting of visualizations allows for greater customization, which can be crucial for effective data communication.

5.2. Neptune.ai

The experience with Neptune.ai highlighted the critical role of experiment tracking within the broader MLOps framework. Tracking our experiments allowed us to document our analysis comprehensively, providing a clear record of the methodologies and results. This documentation fosters reproducibility, a key principle in machine learning, enabling us to revisit or share our findings confidently.

Neptune.ai's simple setup and intuitive interface were particularly advantageous. The platform required minimal configuration, which allowed us to focus on learning how to track our experiments rather than getting caught up in technical details. This ease of use helped build confidence in using the tool effectively.

This project also showed the significance of logging visualizations and results. By storing our visual outputs in Neptune.ai, we ensured that we could revisit our analyses and track progress over time. This practice not only aids in maintaining consistency but also enhances collaboration, as team members can easily access and understand previous work without needing extensive explanations.

5.3. Key Learning Points from the Project

- **The simplicity of `create_report()`:** One of the most significant learning points was how easy it was to generate a comprehensive report using the `create_report()` function. For beginners, this feature greatly boosts productivity by eliminating the need for multiple steps involved in traditional EDA. Instead of writing extensive code, we could get a complete overview of our data in seconds.
- **Using Seaborn for Customized Visualizations:** In addition to DataPrep, we employed Seaborn to create more customized visualizations. By utilizing Seaborn's flexible syntax, we could tailor our plots with specific requirements to try out.
- **Limitations in Customization:** While DataPrep simplifies many tasks, it lacks flexibility for highly tailored visualizations or detailed control over analysis compared to more hands-on libraries like Seaborn and Matplotlib.
- **Enhanced Experiment Management:** Neptune.ai provided a structured framework for organizing and tracking various experiments, which simplified the management of parameters and metrics. This allowed for better oversight of the data analysis and model training processes.
- **Team Collaboration:** The platform enabled improved communication and collaboration among team members. By sharing experiment results and visualizations in a centralized location, all team members had access to the same information, which reduced misunderstandings and improved efficiency.
- **File Uploading:** The functionality to upload and store visualizations and reports in Neptune.ai ensured that all relevant files were organized and easily accessible. This streamlined the workflow and eliminated the need for external storage solutions.
- **Version Control and Comparison:** Neptune.ai's version control features allowed for the effective tracking of different model iterations. This capability made it easier to compare performance metrics across versions, enabling informed decision-making regarding model selection.

6. Conclusion

In this project, utilizing DataPrep and Neptune.ai for Exploratory Data Analysis (EDA) and experiment tracking provided us with valuable hands-on experience in MLOps tools. DataPrep enabled us to automate the EDA process, generating insightful reports and visualizations efficiently. The simplicity of the library allowed us to focus on interpreting the data rather than getting overwhelmed by complex coding tasks. This streamlined approach was particularly beneficial for beginners, as it removed some barriers to entry associated with data analysis.

On the other hand, Neptune.ai proved to be an effective platform for tracking our project version. By facilitating the organization and documentation of our analyses, Neptune.ai enhanced our ability to reproduce results and collaborate effectively. The straightforward setup and minimal configuration required allowed us to engage with MLOps concepts without feeling intimidated, further reinforcing our understanding of experiment tracking in machine learning workflows.

Takeaways:

- **For DataPrep:** The primary advantage lies in its automation capabilities and user-friendly interface, which make EDA accessible and efficient for beginners. This automation helps save time and reduce the learning curve associated with traditional data analysis methods.
- **For Neptune.ai:** The platform excels in providing efficient tracking and promoting reproducibility. Its intuitive design allows users to log their visualizations and results easily, reinforcing the importance of maintaining clear records in the data analysis process.

In conclusion, both DataPrep and Neptune.ai serve as excellent starting points for beginners venturing into the field of MLOps. These tools not only simplify workflows but also lay the foundation for a deeper understanding of more complex machine-learning practices. As we continue our journey in data science, we recognize the importance of leveraging such tools to enhance our productivity and effectiveness in managing machine learning projects.

7. Project Attachments

- Neptune.ai project link:
<https://app.neptune.ai/o/sbonga4/-/projects>
- A downloaded copy of the report created after running the `create_report()` function is attached with the submission.
File name: `titanic_eda_report.html`
- Google Colab code file link:
https://colab.research.google.com/drive/1FiKtbnVKK0brwSqCeAhXRcdgdUL2MN7Q?authuser=0#scrollTo=m9_d2D0V7HmP
- GitHub Link:
<https://github.com/sanjana-bongale/mlops>