

End-to-end Machine Learning pipeline

Author: Sergio Vázquez

Abstract	3
Introduction	4
Context for the pipeline development	5
Architecture	6
Preprocessing	7
Classes	8
PreprocessingPipeline	8
FeatureSelector	8
ClassBalancer	9
DataTransformer	9
Scaler	9
Encoder	10
DataCleaner	10
FeatureCombiner	10
FeatureExtractor	10
Train, Test and Output	11
Classes	12
ModelPipeline	12
Train	12
EvaluationPipeline	12
EvaluateModel	12
OutputModule	13
TrainTest	13
BootstrapPoint632	13
Model	14
Model Training	14
Model Evaluation	15
Usage Guidelines	17
Main usage	17
Implementation of preprocessing methods	17
Implementation of a new Machine Learning algorithm	18
Conclusion	19

Abstract

This comprehensive document aims to elucidate the intricacies of an end-to-end Machine Learning pipeline, shedding light on its practical applications and core objectives, while offering a step-by-step guide on the operational intricacies that underpin its functionality.

The primary mission of this pipeline is to deliver a seamless Machine Learning service, empowering users to take raw, unprocessed dataframes and effortlessly obtain predictions based on their specified target variable. Currently, this pipeline is tailored to handle dichotomous variables, offering a versatile tool for binary classification tasks, making it an invaluable asset for a wide range of applications.

The Machine Learning pipeline outlined here is more than just a technical framework; it represents a bridge between complex data and actionable insights. Whether you're delving into predictive analytics, customer behavior analysis, or fraud detection, this pipeline provides a structured and user-friendly path to harness the power of Machine Learning, even if you're not a seasoned data scientist.

In this document, we will explore the key components of the pipeline, from data preprocessing and feature engineering to model selection and training, followed by deployment and usage. We'll highlight the importance of data quality, model interpretability, and performance evaluation, ensuring that users have a comprehensive understanding of the entire process.

Introduction

The motivation behind this project is driven by a compelling need for an advanced Machine Learning pipeline that goes beyond the limitations of a traditional Jupyter Notebook-based approach. Several key objectives underscore this endeavor, each aimed at enhancing the efficiency, flexibility, and collaboration potential in the world of data science and predictive modeling.

1. **Scalability:** The project is motivated by the necessity to harness new technologies and adapt to evolving data requirements without having to rewrite the entire codebase. By designing a flexible and extensible pipeline, we ensure that our Machine Learning system can readily incorporate novel algorithms, data sources, and processing methods, making it future-proof and adaptable to the ever-changing landscape of data science.
2. **Reproducibility & Version Control:** One of the project's fundamental objectives is to enable robust experimentation and ensure that all changes to the pipeline are well-documented and version-controlled. This commitment to reproducibility is essential for tracking and understanding the evolution of the models and their performance, fostering trust in the results and simplifying the debugging process.
3. **Modularization:** Modularization is pivotal to the project's design, allowing for the replacement or enhancement of individual components without affecting the overall functionality. This design principle promotes code maintainability and eases the process of swapping out parts of the pipeline as the need arises, thereby reducing the risk of system-wide disruptions during updates or improvements.
4. **Performance Optimization:** The project aspires to leverage parallel programming and resource optimization techniques to maximize the efficiency of the Machine Learning pipeline. By distributing tasks across multiple resources and optimizing resource allocation, the system aims to significantly enhance performance, reducing processing times and resource overhead.
5. **Collaboration:** Collaborative data science is at the forefront of our objectives. Enabling multiple team members to contribute to different parts of the pipeline ensures that diverse expertise and insights can be integrated seamlessly. This promotes synergy among team members and fosters the development of a more powerful and comprehensive Machine Learning solution.

In summary, this project is a response to the limitations of traditional Jupyter Notebook-based approaches, aiming to harness the advantages of scalability, reproducibility, modularization, performance optimization, and collaboration.

Context for the pipeline development

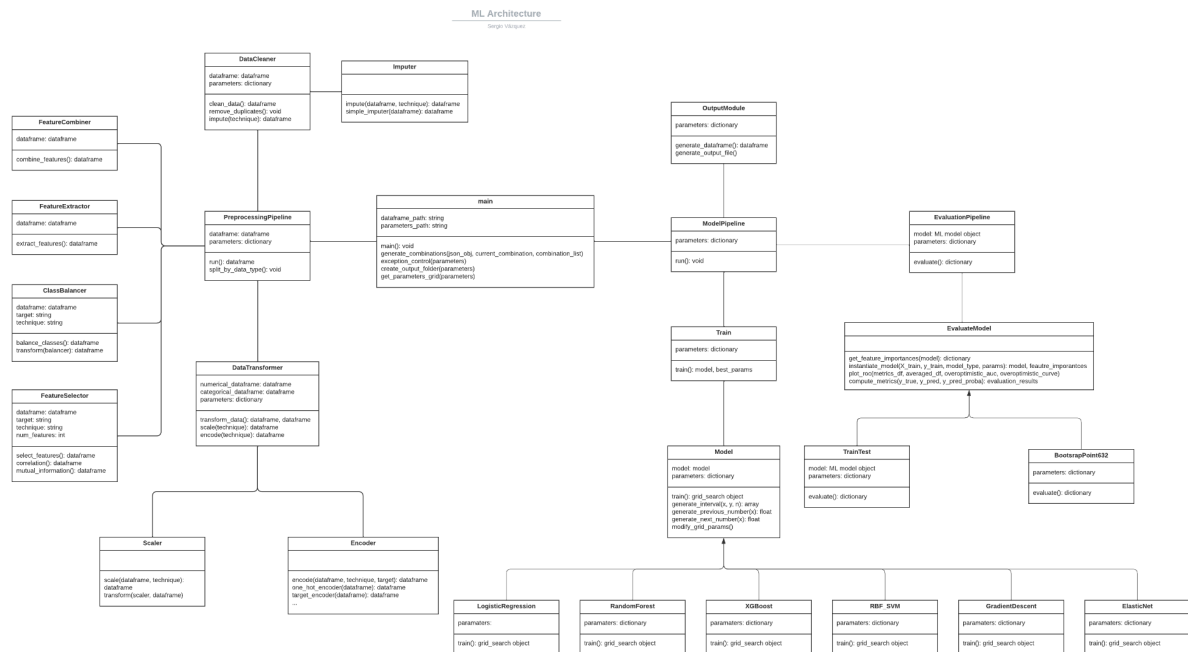
The development of this Machine Learning pipeline emerged from a pressing need within the realm of oncology and precision medicine. Specifically, it was born out of the requirement to construct a predictive model for the efficacy of Immune Checkpoint Inhibitors (ICI) in the context of bladder cancer (BLCA) treatment. This project is an embodiment of the relentless pursuit of improved therapeutic outcomes and personalized medical interventions.

Bladder cancer, one of the most common cancers, presents a unique challenge due to its inherent complexity and the variable response of patients to Immune Checkpoint Inhibitor treatments. These therapies hold great promise but necessitate a highly tailored and precise approach to maximize their effectiveness. This is where the Machine Learning pipeline comes into play, serving as a versatile and indispensable tool in our quest to enhance patient care and treatment success.

At the heart of this endeavor lies the fusion of data from multiple sources. The ability to aggregate and harmonize diverse data sets is pivotal to the success of our predictive model. Our objective is to create a predictive model that can discern the subtleties and intricacies of each patient's condition and deliver a prognosis on the likelihood of a positive response to Immune Checkpoint Inhibitor therapy.

What sets this pipeline apart is its adaptability and reusability. In our relentless pursuit of improved predictive accuracy, we recognized the need for a tool that can seamlessly accommodate changes and innovations. This includes the addition of new variables, the incorporation of cutting-edge techniques, and the ability to experiment with various approaches without the need for extensive code overhauls. This flexibility empowers us to stay at the forefront of research and innovation, ensuring that our predictive model remains agile and effective in the dynamic landscape of medical science.

In summary, the development of this Machine Learning pipeline is more than a technical advancement; it represents a commitment to advancing patient care and addressing the evolving challenges in the field of oncology. It embodies the spirit of collaboration, innovation, and adaptability, all in service of improving the lives of those affected by bladder cancer and bringing us closer to the realization of more precise and effective cancer treatments.



Architecture of the end-to-end ML pipeline

Before delving into the intricate details of each architectural component, the diagram presented above provides an overview of the pipeline's global structure. Positioned on the right side, we find the train/test and output module, while on the left, the preprocessing module is distinctly visible. These two fundamental facets are seamlessly interconnected through the main class, which plays a pivotal role in orchestrating the entire program.

The primary purpose of this central main class is to serve as the conductor of this well-orchestrated symphony of data processing and machine learning. It receives two essential parameters upon execution: the first parameter being the path to a dataframe, and the second parameter designating a path to a comprehensive parameters file. This parameters file functions as the cornerstone of the pipeline, encapsulating all the crucial settings and configurations required to run the pipeline successfully.

This configuration file encompasses a wide array of parameters, covering key aspects of the pipeline's functionality. It allows us to define vital components, such as the choice of machine learning algorithm, the selection of data scaling or balancing techniques, and the destination path for storing the final output. Furthermore, the parameters file is designed to accommodate multiple parameter values for certain fields. This flexibility empowers each run of the pipeline to explore and experiment with diverse combinations of parameters, an invaluable capability in optimizing the predictive model.

Remarkably, the main class assumes the responsibility of generating all possible combinations of parameters, orchestrating multiple runs of the pipeline, and skillfully concatenating the results. This automated and systematic approach ensures that no stone is left unturned in the pursuit of fine-tuning and enhancing the predictive model's performance.

In essence, this robust and dynamic architecture encapsulates a comprehensive and flexible approach to predictive modeling. It not only streamlines the process but also enables the exploration of a multitude of parameter combinations, allowing us to uncover the most effective strategies for predicting responses to Immune Checkpoint Inhibitor therapy in bladder cancer patients. This meticulous and adaptable framework is the cornerstone of our endeavor to provide personalized, data-driven insights in the realm of oncology and healthcare.

Preprocessing

Stepping into the initial segment of the pipeline, we find a crucial phase dedicated to the essential preprocessing and feature engineering procedures. This cornerstone of the pipeline is meticulously designed to pave the way for optimal model performance. Within this preprocessing segment, a multifaceted array of operations is at our disposal. These operations, which can be seamlessly tailored to the unique needs of the dataset, encompass seven distinct steps:

1. **Feature Extraction:** This step is all about distilling the essence of the data, extracting meaningful features that carry the most predictive power for the model. Various techniques can be employed here to identify and extract the most informative attributes.
2. **Feature Combination:** Combining or engineering new features from existing ones can unlock latent insights in the data. This step empowers the pipeline to derive new variables that enhance the model's predictive capacity.
3. **Data Cleaning:** Robust data cleaning is an indispensable part of the process. It encompasses duplicate removal to ensure data integrity and imputation for handling missing values, thereby preparing a pristine dataset for modeling.
4. **Balancing of Classes:** Dealing with imbalanced class distributions is often a crucial consideration in predictive modeling. Techniques for class balancing can be applied to ensure that the model is not skewed by the predominance of one class over another.
5. **Feature Selection:** Not all features are created equally. This step allows for the selection of the most relevant features, streamlining the model and improving its efficiency and interpretability.
6. **Scaling:** Ensuring that features are on a consistent scale is essential for many machine learning algorithms. Scaling techniques standardize the range of features, preventing certain attributes from dominating the model.
7. **Encoding:** Categorical data often requires encoding into numerical form for the model to process it effectively. Encoding techniques are applied to facilitate the integration of these non-numeric variables into the predictive model.

Each of these steps is characterized by flexibility and adaptability, as they can house various techniques and approaches tailored to the specific dataset and project objectives. Furthermore, the pipeline offers the versatility to enable or disable each of these steps, granting complete control over the preprocessing process. This feature proves invaluable when working with preprocessed data or when specific preprocessing steps are deemed unnecessary for a particular analysis.

The central, orchestrating class manages the intricacies of each step, expertly invoking the various classes that handle these processes. It coordinates the workflow, ensuring that each operation is executed seamlessly, while also collecting and synthesizing the results of each step. This central hub is the linchpin that makes the entire preprocessing and feature engineering phase a well-organized and streamlined process, ultimately setting the stage for the subsequent stages of model development and evaluation.

Classes

PreprocessingPipeline

Description: Main class of this part of the architecture, in charge of performing all the preprocessing steps.

Attributes:

- **Dataframe:** Original dataframe to be transformed across the pipeline.
- **Parameters:** parameters dictionary that contain all the information related to the process.

Methods:

- **__init__:** Initializes a new instance of the Pipeline.
- **split_by_data_type:** splits the data in categorical and numerical.
- **run:** main method of the pipeline, it invokes the different classes related to the preprocessing part for finally returning the preprocessed dataframe.

FeatureSelector

Description: This module is in charge of selecting those variables that will be more useful to our predictive model.

Attributes:

- **dataframe:** data.
- **target:** name of the target variable.
- **technique:** technique to be used on the feature selection process.
- **num_features:** number of features to be selected.

Methods:

- **__init__:** Initializes a new instance of the FeatureSelector class.

- `select_features`: Invokes the appropriate feature selection method according to the selected technique.
- `correlation`: implementation of correlation feature selection technique.
- `mutual_information`: implementation of the mutual information technique.

ClassBalancer

Description: Balances the data according to the target variable that we are studying.

Attributes:

- `dataframe`: data.
- `target`: name of the target variable.
- `parameters`: dictionary containing all the needed information.

Methods:

- `__init__`: Initializes a new instance of the ClassBalancer class.
- `balance_classes`: Invokes the appropriate balancing method according to the selected technique.
- `transform`: method in charge of transforming the data.

DataTransformer

Description: This class deals with all the processes related to data transformation, including scaling and encoding.

Attributes:

- `numerical_dataframe`: part of the main dataframe that is numerical.
- `categorical_dataframe`: part of the main dataframe that is categorical.
- `parameters`: parameters to be used.

Methods:

- `__init__`: Initialize a new instance of DataTransformer.
- `transform_data`: this is the main function of the class, it deals with all the transformations by using the parameters object.
- `scale`: auxiliary method used for invoking the scaler.
- `encode`: auxiliar method used for invoking the encoder.

Scaler

Description: Class that is in charge of scaling the original data.

Attributes: -

Methods:

- `scale`: Scales the dataframe using a particular technique.
- `transform`: transforms the dataframe.

Encoder

Description: Class in charge of encoding categorical variables into a format that will be readable for the ML algorithm.

Attributes: -

Methods:

- Encode: invokes the encoding technique that was set via parameters.
- One_hot_encoder: implementation of one hot encoding.
- Target_encoder: implementation of target encoding.

DataCleaner

Description: Class dealing with standardization of columns, duplicate removal, data formatting, imputation and outlier removal.

Attributes:

- dataframe: data.
- parameters: parameters dictionary that contains all the information related to the process.

Methods:

- __init__: initializes a new instance of the DataCleaner class.
- clean_data: This main method will perform the data cleaning procedure.
- remove_duplicates: removes all the duplicated rows from the dataframe.
- impute: calls the imputer in order to perform the imputation of the dataframe.

FeatureCombiner

Description: This module is meant to capture more complex relationships between variables, including multiplication between them or any other linear combination of them.

Attributes: -

Methods: -

FeatureExtractor

Description: This module will extract new features from the existing ones in order to add more valuable information to the model by using the already reported features.

Attributes: -

Methods: -

Train, Test and Output

Within the architecture, the second segment is dedicated to the critical phases of training, testing, and the ultimate output of results. This multifaceted module plays a pivotal role in the pipeline's journey toward constructing a reliable predictive model. It can be further dissected into three distinct but interconnected parts:

- **Training Procedures:** This initial part of the module is the crucible where the model takes shape. It encompasses the entire training process, a cornerstone of machine learning. Here, a series of pivotal activities unfold. Hyperparameter tuning is meticulously executed, fine-tuning the model's configurations to extract the utmost performance. Concurrently, feature importances are derived, shedding light on which attributes exert the most influence on the model's predictions. The culmination of this phase is the trained model itself, a dynamic entity primed to make insightful predictions.
- **Testing and Evaluation:** The next segment, the testing module, assumes the responsibility of subjecting the trained model to rigorous evaluation. It operates with different strategies, the first one is reserving a part of the data for testing, enabling the computation of a comprehensive array of machine learning performance metrics, the second one includes bootstrapping the initial dataset in order to perform internal validation. These metrics, including but not limited to AUC (Area Under the Receiver Operating Characteristic Curve), F1 score, and recall, provide a holistic understanding of the model's predictive capabilities. The results of this phase are a litmus test of the model's proficiency, offering insights into its real-world performance.
- **Output and Reporting:** The output module, the final act of this segment, is tasked with encapsulating the comprehensive scope of the pipeline's execution. It diligently gathers all relevant information, encompassing the parameters and hyperparameters employed, feature importances, and the full spectrum of performance metrics. This data is meticulously compiled and then consolidated into a structured format, usually a CSV file. The resultant file is conveniently saved to a specified path, serving as a comprehensive record of the pipeline's execution, insights, and outcomes. Furthermore, a report file is created that will compress all the important information into a txt file.

In sum, this module harmonizes the key aspects of model development and evaluation, culminating in a comprehensive understanding of the model's predictive prowess. It is a testament to the rigor and precision that underpins the pipeline, providing both data scientists and stakeholders with an invaluable resource for informed decision-making and further refinement of the model. This trifecta of training, testing, and reporting forms the backbone of the predictive modeling process, making it an integral part of the journey toward developing a robust Immune Checkpoint Inhibitor therapy response predictor for bladder cancer patients.

Classes

ModelPipeline

Description: Pipeline in charge of running all the train/test/evaluation procedures.

Attributes:

- parameters: parameters dictionary containing all the needed information.

Methods:

- __init__: Constructor used for creating the ModelPipeline
- run: main method of the class, in charge of creating the different objects involved in the procedure.

Train

Description: Initializes a new instance of the Train class.

Attributes:

- parameters: parameters to be used.

Methods:

- __init__: Constructor of the Train class.
- get_feature_importances: method used for obtaining the feature importances after training the model.
- train: performs the training step which varies depending on the type of model that is selected.

EvaluationPipeline

Description: Pipeline that is in charge of performing the evaluation of the machine learning model.

Attributes:

- model: machine learning model.
- parameters: dictionary that contains all the parameters needed for the execution.

Methods:

- evaluate: core method of the class, in charge of invoking the different classes that will perform the evaluation.

EvaluateModel

Description: Superclass for the different evaluation methods, contains all the necessary code for them.

Attributes: -

Methods:

- `get_feature_importances`: method used for obtaining the feature importances after training the model.
- `instantiate_model`: creates a particular model, fits it and returns the fitted model along with its feature importances.
- `plot_roc`: method for plotting and saving the roc curve.
- `compute_metrics`: method that contains all the code for the metrics.

OutputModule

Description: Class in charge of allocating all the information of the run inside a csv file.

Attributes:

- `parameters`: parameters dictionary that will contain all the information of the run, including the training and testing results.

Methods:

- `__init__`: Initializes a new instance of the Output Module.
- `generate_dataframe`: Method used for generating the dataframe that will store the output of both pipelines.
- `generate_output_file`: Method that will generate the report encapsulating all the information about the pipeline run.

TrainTest

Description: Class in charge of performing the train/test evaluation technique.

Attributes:

- `parameters`: dictionary of parameters.
- `model`: machine learning model.

Methods:

- `evaluate`: method that performs the evaluation by following the train/test technique.

BootstrapPoint632

Description: Class in charge of performing bootstrap .632+.

Attributes:

- `parameters`: dictionary of parameters.

Methods:

- `evaluate`: method that performs the evaluation by Bootstrapping .632+.

Model

Description: Superclass created for englobe all the Machine Learning models, which represents the concept of Model which is meant to store all the methods commonly used for training models.

Attributes:

- x: Part of the data that is going to be used for training the model.
- y: Target vector.
- model: Machine learning model to be used.
- param_grid: hyperparameters grid to be used during the grid search.

Methods:

- `__init__`: Constructor class of the Model superclass.
- Train: Method that recursively trains the model modifying the hyperparameters dynamically until the current score do not improve the old one.
- Generate_interval: generates an interval between x and y of size n.
- Generate_previous_number: Used to generate a previous number for the hyperparameter fine tuning when the optimal value is the lowest one in the list.
- Generate_next_number: Used to generate the next number for the hyperparameter fine tuning when the optimal value is the highest one in the list.
- Modify_grid_params: method that modifies all the numerical list of hyperparameters dynamically. It iterates through all the hyperparameters and taking into account the position of the optimal value, it generates a new range of values that will be used in order to optimize even more the hyperparameters.

Model Training

The linchpin of the training process within this architecture is the central Train class, which is summoned into action from the ModelPipeline. Once the main method is invoked, the pipeline orchestrates the entire training phase, tailoring its operations according to the chosen Machine Learning algorithm.

The beauty of this design is the hierarchical structure that underlies the creation of specific classes for each algorithm. If a particular algorithm is to be employed, all that's needed is the creation of a corresponding class dedicated to that algorithm, with each of these classes extending the Model superclass, which acts as the repository for the core code. After crafting this subclass, the only remaining step is to instantiate it and pass it to the superclass.

The train method within the Model superclass operates as the maestro of the training process. It initiates a Cross-Validated Grid Search, a systematic approach to hyperparameter optimization. This search not only fits the data but also tirelessly seeks the best set of hyperparameters to maximize model performance.

The best thing about this method lies in its adaptability. Following each round of the grid search, the grid parameters are cleverly adjusted based on the optimal parameters identified in the previous run. New intervals are generated, allowing the grid search to refine the algorithm further. This iterative process continues until the performance score, which in this case is measured by the ROC AUC (Receiver Operating Characteristic Area Under the Curve), no longer shows improvement.

Once this intricate process concludes, the trained model emerges as the final result. Notably, this isn't the end of the line; feature importances are meticulously extracted from the model, shedding light on the attributes that wield the most influence over predictions. Furthermore, the best hyperparameters are also recorded, providing a clear snapshot of the optimized configuration for the specific algorithm.

Model Evaluation

The model evaluation process is a critical phase, given that the sample size is insufficient to rely solely on the traditional train/test approach. Therefore, the pipeline employs various evaluation techniques, including train/test, and bootstrap .632+.

In the train/test approach, data is separated after all preprocessing steps. The model is trained on a portion of this data and tested on a distinct test cohort that the model has never encountered. This process is repeated a specified number of times to mitigate the impact of chance and provide robust results.

The bootstrap .632+ follows a similar approach, but the model is trained on a bootstrap sample and tested on Out-Of-Bag (OOB) samples, which are samples not included in the training dataset (as Bootstrap involves duplicating samples). The .632+ version calculates final metrics through a weighted computation, considering metrics obtained from testing on both the same data used for training and the OOB samples.

Finally, it is worth mentioning that all different evaluation techniques make use of common code for computing the different metrics, these metrics encompass:

- Accuracy: A fundamental metric that gauges the proportion of correctly predicted instances out of the total number of instances in the test set.
- Precision: This metric highlights the model's ability to make positive predictions accurately, measuring the ratio of true positives to the sum of true positives and false positives.
- Recall: Also known as sensitivity or true positive rate, recall quantifies the model's capability to identify all actual positive cases by calculating the ratio of true positives to the sum of true positives and false negatives.
- F1 Score: The F1 score harmonizes precision and recall, providing a balanced measure of a model's accuracy that is especially useful when dealing with imbalanced datasets.

- AUC (Area Under the Curve): The AUC assesses the model's ability to distinguish between positive and negative instances, as depicted by the Receiver Operating Characteristic (ROC) curve. A high AUC value indicates strong predictive performance.
- Confusion Matrix: A visual representation of the model's classification performance, the confusion matrix provides insights into true positives, true negatives, false positives, and false negatives, aiding in the understanding of prediction behavior.
- Brier score: Is a measure used for assessing the accuracy of probabilistic predictions, a lower Brier score indicates that the model's predicted probabilities are more accurately reflecting the true probabilities of an outcome occurring.

These metrics collectively offer a comprehensive perspective on how well the model is performing in its predictive capabilities. They help assess the model's ability to correctly classify instances, its robustness in capturing true positives, and its overall predictive power in distinguishing between the two classes.

The results of these evaluations provide a holistic understanding of the model's accuracy, precision, and recall, along with its ability to balance trade-offs between precision and recall via the F1 score. Furthermore, the AUC and confusion matrix offer insights into the model's capacity to discriminate between classes and its classification tendencies.

In essence, the evaluation module serves as the litmus test for the model's real-world performance, furnishing vital insights that guide decision-making and inform any necessary adjustments or refinements to the predictive model.

Usage Guidelines

Main usage

There is no particular way of executing the program, as said in this document, the only two parameters that are needed are the path to the dataframe and the path to the parameters file, the parameters file structure will be appended at the end of this document. Once you have this two paths, the only thing you need to do is:

```
py mlpipeline.py path1 path2
```

Implementation of preprocessing methods

In a true testament to the flexibility and extensibility of this Machine Learning pipeline, the process of implementing new techniques for any of the preprocessing steps is remarkably straightforward. The structure and architecture of this pipeline have been thoughtfully designed to make such enhancements seamless and efficient.

When an individual desires to introduce a new method or technique within a specific preprocessing step, there's no need to embark on an extensive overhaul of the entire pipeline. Instead, one simply navigates to the particular class associated with the desired preprocessing step. Within this class, the proposed structure offers an elegant and intuitive approach to accommodating new techniques.

This structure typically revolves around a tree of conditional statements, typically "if" statements, which dynamically respond to the selected technique. Each "if" block corresponds to a different technique, executing specific code tailored to the chosen method. The process of introducing a new technique becomes as straightforward as adding another "if" block, specifying the technique to be implemented, and providing the corresponding code necessary for this technique to seamlessly integrate into the pipeline.

Finally, it is **important** that the new preprocessing method is added to the handling exceptions part of the pipeline (main.py), as it only accepts those values that are already implemented.

A key point of emphasis is that the core functionality of the main pipeline remains untouched. The integrity and robustness of the pipeline are maintained by confining the code extensions to the specific class where the new technique is being added. This clear separation of concerns ensures that the pipeline's overall structure remains unaltered, preserving its reliability and ease of maintenance.

Implementation of a new Machine Learning algorithm

Expanding the repertoire of machine learning algorithms within this pipeline is as streamlined and efficient as the incorporation of new preprocessing techniques, if not more so. As elucidated in the architecture section, the process is designed for simplicity and extensibility, underscoring the pipeline's adaptability.

The key to this simplicity lies in the Model superclass. This superclass serves as the repository for all the essential code required to execute the training phase of the pipeline. Therefore, if there's a desire to introduce a new algorithm to the pipeline, the process is remarkably straightforward:

1. **Class Extension:** A new class is created, which should extend the Model superclass. This straightforward act establishes the foundation for introducing the new algorithm to the pipeline.
2. **Constructor Implementation:** Within this new class, two methods need to be implemented. The first is the constructor, which is remarkably concise. It simply calls the constructor of the superclass, passing as a parameter the particular algorithm used. In many cases, this corresponds to the sklearn implementation of the algorithm.
3. **Parameter Grid Inclusion:** This parameter grid should be included as a class parameter within the new class, containing the particular hyperparameters associated with this algorithm. This grid defines the space to be explored during hyperparameter tuning.
4. **Add the keyword for the new ML model to the handling exception part of the pipeline,** so it does not throw an exception when using this model (main.py).
5. **Add the model invocation in `instantiate_model:EvaluateModel`.**
6. **In case the model needs any particular way of extracting feature importances,** this has to be implemented on `instantiate_model:EvaluateModel`, the implementation is quite straightforward and one just needs to follow the if/else structure.
7. **Training Method Invocation:** Finally, the new class should invoke the train method that is implemented in the superclass. This method handles the entire training process, including hyperparameter tuning and feature importance extraction. By simply invoking this method, the new algorithm is seamlessly integrated into the pipeline's training phase.

Conclusion

In summary, the development of this Machine Learning pipeline represents an approach to the complex task of predicting responses to Immune Checkpoint Inhibitor therapy in bladder cancer patients. The motivation behind this endeavor is grounded in the urgent need for precise, data-driven solutions in the field of oncology. With a strong focus on scalability, reproducibility, modularization, performance optimization, and collaboration, this pipeline promises to revolutionize the way predictive models are created and maintained.

The architecture of the pipeline is meticulously designed, with a clear separation of preprocessing, training, and evaluation phases. The preprocessing module offers flexibility through a range of operations, allowing for data preparation, feature engineering, and the handling of class imbalances, all with the ability to enable or disable specific steps based on the project's needs.

The training phase is highly adaptive, employing a hierarchical structure that caters to various Machine Learning algorithms. The grid search technique fine-tunes hyperparameters, ensuring the model's optimal performance, while feature importances and parameter records provide invaluable insights into the model's inner workings.

The evaluation module rigorously tests the model on a separate test dataset, computing a comprehensive set of metrics that assess its accuracy, precision, recall, F1 score, AUC, and even the confusion matrix. These metrics offer a holistic view of the model's predictive capabilities, guiding decisions and refinements.