

Automated Machine Learning Orchestration Using PYCARET and STREAMLIT

C. Thangalatha Legaz¹, S. Anjum², R. V. Pooja², M Sowmiya²
Assistant Professor¹, UG Scholar²

Department of Information Technology,
RAAK College of Engineering and Technology,
Puducherry, India.

thangalatha@gmail.com

Abstract – The convergence of AutoML and data preprocessing represents a significant advancement in machine learning, addressing key challenges in model development. Tree-based Pipeline Optimization Tool (TPOT) streamlines algorithm selection and hyperparameter tuning while efficiently handling preprocessing steps like class imbalances and feature selection. TPOT's capabilities are complemented by PyCaret, offering a suite of tools for automating data preprocessing tasks such as feature engineering and encoding. PyCaret's accessibility and support for 15 algorithms, including advanced techniques, broaden the scope of model experimentation. Together, TPOT and PyCaret simplify the machine learning workflow, making it accessible to both experts and novices. This combination enhances model performance and fosters the ongoing evolution and democratization of machine learning, advancing its application across various domains.

Index Terms –AutoML, Data, Preprocessing, MachineLearning, Hyperparameters, Feature selection, Report generation, Data Visualization

1. INTRODUCTION

The presented paper navigates the intricate relationship between Machine Learning (ML) and AutoML, emphasizing their symbiotic collaboration. ML, relying on vast datasets, employs algorithms to discern patterns, make predictions, and facilitate decision-making. AutoML complements this by automating essential yet tedious tasks, enhancing the efficiency and accessibility within the ML pipeline.

In the pursuit of refining data preprocessing methodologies, the survey proposes a comprehensive exploration of research papers and contributions. This extensive study encompasses a spectrum of papers that delve into innovative solutions bridging the gap between current challenges and the future promises of data preprocessing. The objective is to extract crucial insights by merging the intricacies of data preprocessing with the transformative potential of AutoML, aiming to advance data-driven decision-making in the evolving ML landscape. Highlighted within this survey are various research papers, from "DataAssist" to "REIN," each contributing vital principles to guide this quest. These papers shed light on the complex terrain where challenges such as imbalanced data, hyperparameter optimization nuances, and the need for advanced feature engineering converge, necessitating holistic solutions to bridge the divide between data and model.

As the frontiers of machine learning continue to expand, the principles extracted from these research papers act as guiding beacons, urging the crafting of automated solutions capable of addressing multifaceted challenges. Rooted in this extensive literature survey, the forthcoming architectural overview promises not only innovation but a transformation of the status quo. It aims to provide all-encompassing, end-to-end solutions for data preprocessing and the automation of pivotal tasks. The survey's principles revolve around problem identification and the search for innovative solutions, bridging the gap between present challenges and the promising future of data preprocessing. Each research paper in the survey addresses a specific facet of data preprocessing, collectively contributing to a comprehensive understanding of this critical domain.

2. RELATED WORKS

Runtime Prediction of Machine Learning Algorithms in Automated Systems Parijat Dube; Theodoros Salonidis [1] DataAssist sounds like a promising addition to the landscape of automated machine learning (AutoML) tools. By focusing on data preparation and cleaning, it addresses a crucial aspect of the machine learning workflow that is often overlooked by existing tools. The key features of DataAssist seem to align well with the needs of data scientists and analysts, particularly in industries where data quality is paramount, such as economics, business, and forecasting. By providing functionalities for exploratory data analysis, visualization, anomaly detection, and data preprocessing, DataAssist streamlines the process of preparing data for modeling, potentially saving significant time and effort for practitioners. Moreover, the ability to export cleaned and preprocessed datasets for integration with other AutoML tools or user-specified models enhances its versatility and interoperability within existing workflows. This flexibility is crucial for accommodating different preferences and requirements in data analysis pipelines. Overall, DataAssist appears to fill a significant gap in the existing landscape AutoML tools by prioritizing data-centric tasks and offering comprehensive support for data preparation and cleaning. Its potential to save over 50% of the time typically spent on these tasks underscores its value proposition for practitioners across various domains.

Suraj Juddoo investigates data repair steps for EHR Big Data.[2] This paper addresses a significant challenge with relation to big data systems, particularly focusing on Electronic Health Records (EHR). The emphasis on optimizing data quality methodologies aligns well with the growing importance of leveraging high-quality data for meaningful insights, especially in sensitive domains like healthcare. The recognition of the data repair stage as a critical component of the The data quality life cycle involves crucial, as addressing dirty data is often a complex and resource-intensive task. The acknowledgment an ignorance of how well-performing data restoration tools and algorithms work in the context of big data is an important observation, highlighting the need for specialized solutions in this domain. The systematic examination of data repair techniques and tools, then an experiment-based evaluation, is a robust methodology for gaining insights into their effectiveness. The comparison with a prototype built from previous study results adds a practical dimension to the evaluation. The finding that For Big Data, no algorithm or tool was found to be exceptionally sufficient emphasizes the challenges in this domain. However, identifying some algorithms and tools as marginally better than others provides valuable insights for potential improvements. The recommendations for enhancing data repair algorithms and tools for Big Data represent a valuable contribution to the field, guiding future research and development efforts.

Assessing the performance of AutoML algorithms using a set of simulated classification tasks Henrique Pedro Ribeiroa and Patryk Orzechowski [3] This paper explores the growing landscape of The popularity of machine learning automated (AutoML) programs can be attributed to their great performance and versatility in solving a wide range of issues. The challenge lies in choosing the most suitable AutoML algorithm for a given problem amid the increasing options available. To address this, the study examines the output of four well-known AutoML algorithms using their Diverse and generative ML benchmarking (DIGEN): Auto-Sklearn, H2O AutoML, Auto-Sklearn 2, and Tree-based Pipelines Optimizing Tool (TPOT). Synthetic datasets called DIGEN are used to demonstrate the advantages and disadvantages of popular machine learning algorithms. The outcomes demonstrate how successfully AutoML detects pipelines across datasets. While the majority of AutoML algorithms demonstrated similar performance, subtle differences emerged based on specific datasets and evaluation metrics, providing valuable insights into their comparative effectiveness.

A Whole-System Benchmarking Structure for Data Cleaning Techniques in Machine Learning Pipelines: REIN Christian Hammacher, Harald Schoening, and Mohamed Abdelaal [4] The paper emphasizes the crucial role of machine learning (ML) in daily life and emphasizes how important high-quality data is throughout the ML application lifecycle. It acknowledges the common discrepancies present in real-world tabular data, include inconsistencies, duplication, outliers, missing values, and pattern violations, which often arise during data collection, transfer, storage, or integration. Despite numerous data cleaning methods addressing these issues, the paper points out a gap in considering downstream ML model requirements. To bridge this gap, the work introduces a comprehensive benchmark named REIN1, aiming must carefully evaluate how various ML models are affected by data cleaning techniques. The benchmark addresses key research questions, exploring the necessity and efficacy of data cleaning in ML pipelines. The evaluation involves 38 error detection and repair methods, ranging from simple to advanced. To provide

comprehensive insights, the benchmark employs a broad range of machine learning models that were trained on 14 publicly-accessible datasets that span multiple domains and include both synthetic and realistic error characteristics.

AutoCure: Machine Learning Pipeline Automatic Tabular Data Curation Method Ahmad Schoening, Harald Schoening, and Rashmi Koparde [5] The paper introduces Data curation pipeline AutoCure is innovative and requires no setting designed to address the persistent challenge of data preparation in machine learning applications across domains like autonomous driving, healthcare, and finance. The need for expert knowledge and considerable time investment in navigating the extensive search space for suitable data curation and transformation tools is a recognized hurdle in model development. AutoCure stands out by synthetically enhancing the clean data fraction by combining a data augmentation module with an inventive adjustable ensemble-based error detection technique. Notably, its configuration-free nature streamlines the implementation process, making it accessible for integration using free and open-source resources such as Auto-sklearn, H2O, and TPOT, therefore advancing the general democratization of machine learning.

Reciprocal neural networks for bidirectional mistake detection in databases Holzer and Stockinger, Kurt[6] This paper presents an innovative architecture leveraging bidirectional recurrent neural networks for the purpose of error detection in databases. Through experiments conducted on six distinct datasets, the outcomes demonstrate how well this strategy performs in comparison to cutting-edge mistake detection technologies. Specifically, the average F1-scores across all datasets demonstrate the effectiveness of the proposed architecture. Notably, the system exhibits a lower standard deviation, indicating greater robustness compared to existing methods. An additional advantage is the system's ability to achieve high F1-scores without the need for supplementary data augmentation techniques. This signifies the potential of the introduced bidirectional recurrent neural network architecture as a robust and efficient solution for error detection in diverse database scenarios.

3. PROPOSED METHODOLOGY

A variety of technologies were employed in the study, including information collection, dataset preparation, and model evaluation.

DATA PREPROCESSING MODULE:

The role of cleaning and preparing raw data for analysis is a crucial step in the data science workflow, as it significantly influences the accuracy and reliability of subsequent analyses and machine learning models. This responsibility involves a series of tasks aimed at ensuring the data is presented in a format that is appropriate and consistent for meaningful interpretation. One fundamental aspect is handling missing values, where techniques such as imputation or deletion are employed to address the absence of information. Scaling features is another essential task, particularly when variables are measured on different scales, to prevent certain features from disproportionately influencing the analysis. Additionally, encoding categorical variables is necessary to transform qualitative input into a numerical form that machine learning algorithms can handle. This process helps maintain the integrity of the data and ensures that the chosen analytical techniques can effectively derive insights. Overall, the meticulous cleaning and preparation of raw data form the foundation for robust and reliable data analyses, facilitating informed decision-making in various domains.

AUTOML CORE MODULE:

The central module orchestrating the entire AutoML process serves as the backbone of the automated machine learning workflow, playing a pivotal role in coordinating and managing various tasks. This module integrates sub-modules that collectively contribute to the comprehensive AutoML pipeline, ensuring a streamlined and efficient process. Among these sub-modules, hyperparameter tuning is responsible for optimizing the configuration settings of machine learning models to enhance their performance. Feature engineering involves transforming and selecting features to improve The capacity of the model to identify links and patterns in the data. Model selection, another critical sub-module, aids in choosing the most suitable algorithm or ensemble of algorithms for a given task. By consolidating these sub-modules, the central module makes ensuring the AutoML process is coherent and well-coordinated, with each step contributing to the final result. of automating the model development lifecycle and delivering optimized, high-performing machine learning models.

CATEGORICAL VARIABLE ENCODING STANDARDIZATION MODULE:

The task of ensuring consistent encoding of categorical variables is vital in both regression and classification tasks within the context of machine learning. Categorical variables, representing qualitative data, need to be converted into a numerical representation so that it can work with other algorithms. The responsible module addresses this by employing encoding methods that maintain consistency across tasks. Common techniques include label encoding and one-hot encoding, in which binary columns indicate each category, which provides a distinct number label for every category, and target encoding, where categories are encoded based on the mean of the target variable. By implementing these encoding methods consistently, the module ensures that the machine learning models receive uniform input representations, fostering accuracy and reliability in predictions across both regression and classification scenarios. This consistency is essential for creating robust and interpretable models that can effectively learn patterns from categorical features.

USER INTERFACE (UI) MODULE:

The user-friendly interface serves as the gateway for practitioners to interact seamlessly with the AutoML system. Its primary function is to provide an accessible platform where users can input their data, define relevant parameters, and visualize the results of the automated machine learning process. Through an intuitive design, practitioners can effortlessly upload datasets, specify preferences for hyperparameters or feature engineering, and easily navigate through the system's functionalities. The interface abstracts the complexities of the underlying AutoML algorithms, making it suitable for users of different skill levels. Visualization tools incorporated into the interface enable users to interpret and comprehend the outcomes of the automated processes, fostering a transparent and interactive user experience. Overall, the user-friendly interface enhances the usability of the AutoML system, facilitating effective collaboration between machine learning practitioners and the automated system for streamlined model development.

REPORT GENERATION MODULE:

The deployment module plays a crucial role in facilitating the seamless integration of AutoML-generated models into production environments. Its primary function is to streamline the transition from model development to real-world applications. This module often includes features for model versioning, allowing practitioners to track and manage different iterations of models. Additionally, it addresses scalability concerns, ensuring that the deployed models can handle varying workloads and adapt to changing data volumes. Moreover, monitoring capabilities are integrated to keep track of model performance in real-time, enabling timely interventions if issues arise. By encompassing these functionalities, the deployment module enhances the reliability, scalability, and maintainability of AutoML-generated models in production, ultimately supporting the practical and sustainable application of machine learning solutions.

ARCHITECTURE DIAGRAM:

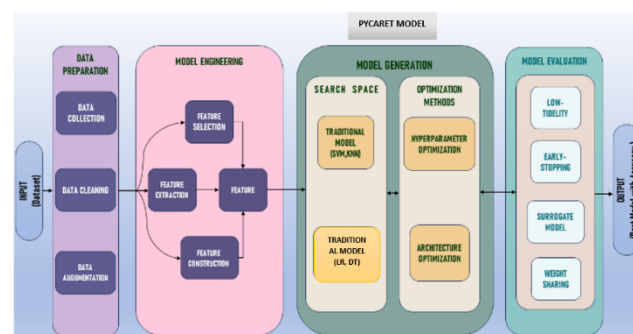


Fig 1: System Architecture

An architecture diagram makes the structure's visual representation available and components of a system or application. It typically includes various elements such as modules, databases, servers, and their interactions. The diagram serves as a high-level overview, illustrating how different parts of the system are connected and work together to achieve the intended functionalities. This visual representation aids in understanding the overall design,

dependencies, and flow of data or processes within the architecture. It is a valuable tool for communication among stakeholders, allowing developers, architects, and other team members to have a shared understanding of the system's structure, helping in decision-making, troubleshooting, and system documentation.

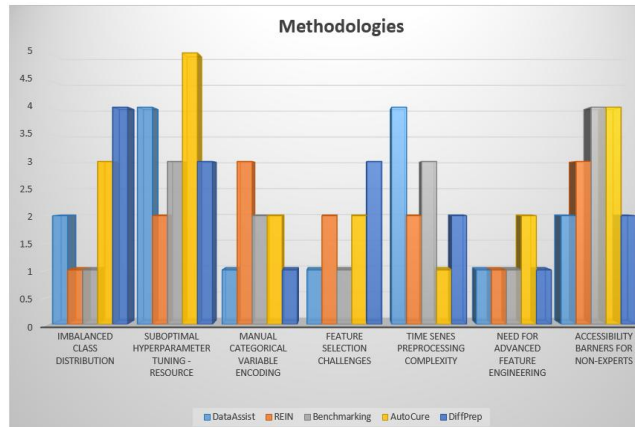
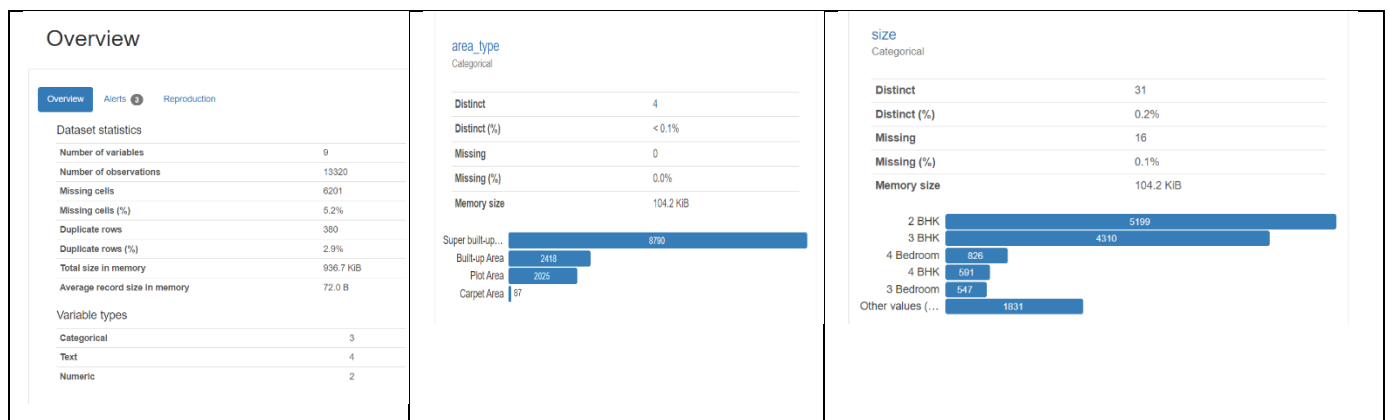


Fig 2: Methodologies

The research paper aims to explore the intricate details of the AutoML system, providing an in-depth analysis of its capabilities, experimental results, and its potential to revolutionize the field of machine learning. With a particular focus on addressing the shortcomings in existing data preprocessing methodologies, the system is positioned as a promising approach to enhancing datasets and subsequently improving findings across diverse domains. The paper likely delves into the system's innovative features, experimental validations, and how it contributes to overcoming challenges in data preprocessing, ultimately paving the way for more effective and efficient machine learning applications. The emphasis on improving datasets suggests a commitment to elevating the overall quality of input data, which is essential to machine learning models' ability to succeed.

4. RESULT AND DISCUSSION

AutoML (Automated Machine Learning) is a broad term that encompasses a variety of techniques and methodologies to automate The procedure for creating ML models. The specific formulas and methods used in AutoML systems can indeed vary based on the tasks being automated.



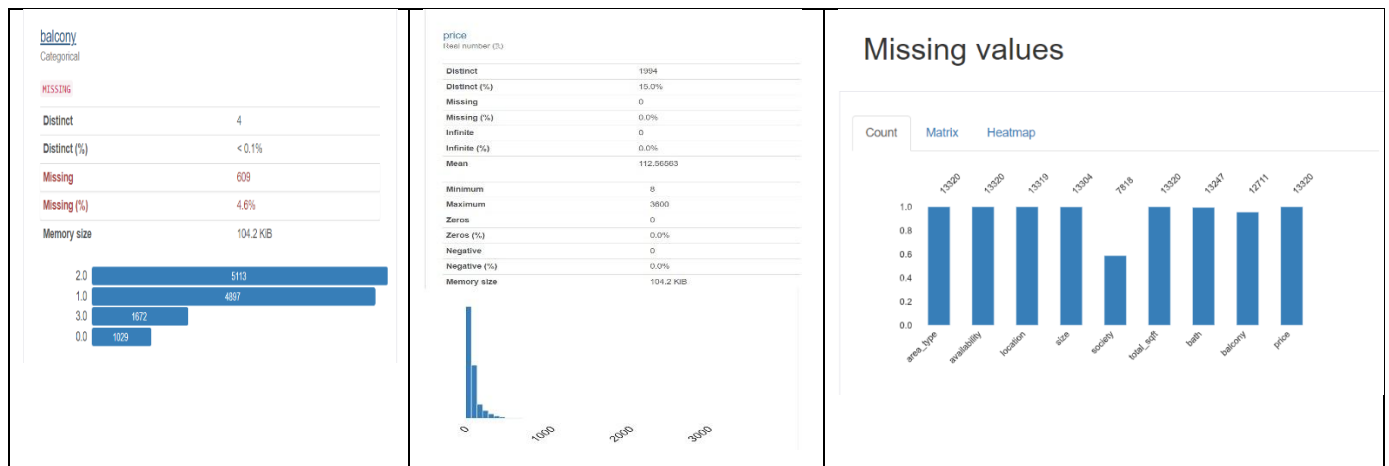


Table 1: Individual Values from the Dataset.

These are some outputs and overview of the given dataset. The given bar graphs are the analysed values of the given dataset which consists of no. of variables, Missing cells, Duplicate rows, Total Size in Memory and Missing Values.

HYPERPARAMETER TUNNING:

Hyperparameter tuning is an essential component in machine learning model optimization, involving the adjustment of external settings that influence the learning process and model performance. These configurations, referred to as hyperparameters, are pre-established before training begins and cannot be learned from the learning set. The goal of hyperparameter tuning is to find the most effective configuration that maximizes a predefined performance metric. Grid search and random search are two commonly employed techniques for this purpose.

$$\text{Best Hyperparameter Value} = \arg \max_{\text{Hyperparameter Values}} \text{Model Performance Metric}$$

Grid search examines a preset set of combinations of hyperparameters methodically, exploring the entire search space. In contrast, random search randomly samples configurations, offering a more stochastic approach. To ensure that the model performs as well as possible on data that hasn't been seen before, both approaches try to find a compromise between generalization and model complexity. A key component of fine-tuning models is hyperparameter tuning for specific tasks, ultimately enhancing their predictive capabilities and robustness.

FEATURE ENGINEERING:

Indeed, feature engineering is a critical aspect of the machine learning model development process, encompassing various techniques to enhance the representation of data and improve model performance. This process involves creating new features, transforming existing ones, or selecting a subset of features to provide the model with more relevant and informative input.

$$\text{New Feature} = \text{Feature}^2$$

Creating new features may involve combining or synthesizing existing features to capture higher-order relationships or patterns in the data. Transformation of features can include normalizing or scaling numerical features, handling missing values, or encoding categorical variables. Additionally, finding and keeping the most essential features is the goal of feature selection while discarding less important ones, reducing dimensionality and potentially mitigating overfitting.

ENSEMBLE METHODS:

AutoML frequently leverages ensemble methods as a powerful strategy to boost overall model performance. Ensemble methods involve combining predictions from multiple individual models, often of diverse architectures or trained with

different subsets of data. The goal is to exploit the complementary strengths of various models, mitigating individual weaknesses and improving overall predictive accuracy. Common ensemble techniques include bagging, boosting, and stacking. Bagging, such as in Random Forests, aggregates predictions from several decision trees that were trained using arbitrary portions of the data improved robustness and decreased overfitting.

$$Ensemble Prediction = \frac{1}{n} \sum_{i=1}^n Model_i (Input Data)$$

Boosting, exemplified by algorithms like AdaBoost or Gradient Boosting, sequentially trains models, with each subsequent model focusing on correcting the errors of its predecessor, leading to increased accuracy. Stacking combines predictions from different models using a meta-model, learning to weigh individual model outputs optimally. Ensemble methods are effective in handling complex relationships within data, increasing model stability, and generalizing well to unseen instances, making them a valuable tool in the AutoML toolkit for achieving superior predictive performance.

MODEL SELECTION:

In AutoML, selecting the best-performing model is pivotal and relies heavily on evaluating various performance metrics. Common metrics include area under the Receiver Operating Characteristic (ROC) curve, accuracy, and F1-score. Accuracy measures the proportion of correctly predicted instances, offering a straightforward assessment of overall correctness. F1-score strikes a balance between recall and precision, making it ideal for applications where the costs of false positives and false negatives differ. The area under the ROC curve quantifies the trade-off between the true positive rate and false positive rate, indicating a model's ability to distinguish between classes. The choice of metric depends on the dataset's nature; accuracy is suitable for balanced datasets, while F1-score is preferable for imbalanced classes. With PyCaret, model creation becomes streamlined, empowering users to leverage a diverse array of algorithms and performance evaluation techniques to select the most suitable model for their specific application.

$$Best Mod = \arg \max_{Models} Model Performance Metric$$

AutoML systems often perform a systematic search over hyperparameter configurations, and The model that performs the best according to the selected metric is the one that gets deployed. These metrics guide the AutoML process, ensuring the chosen model aligns with the specific objectives and requirements of the given machine learning task

	Model	MAE	MSE	RMSE	R2	RMSLE	MAPE	TT (Sec)
gbr	Gradient Boosting Regressor	0.7224	1.0575	1.028	0.0924	0.4842	0.1577	0.177
lightgbm	Light Gradient Boosting Machine	0.7035	1.0753	1.0364	0.0773	0.4813	0.155	0.144
rf	Random Forest Regressor	0.6797	1.0811	1.0392	0.0712	0.4773	0.1527	0.409
lar	Least Angle Regression	0.7809	1.1385	1.0664	0.0237	0.4977	0.1736	0.015
lr	Linear Regression	0.7825	1.1396	1.0669	0.0228	0.4982	0.1739	0.772
ridge	Ridge Regression	0.7826	1.1396	1.0669	0.0228	0.4982	0.1739	0.01
br	Bayesian Ridge	0.7848	1.1402	1.0672	0.0223	0.4984	0.1747	0.012
et	Extra Trees Regressor	0.6681	1.1443	1.0691	0.0169	0.4867	0.1503	0.208
en	Elastic Net	0.7975	1.15	1.0718	0.0139	0.5002	0.1793	0.01
lasso	Lasso Regression	0.8003	1.1524	1.0729	0.0119	0.5008	0.18	0.012
llar	Lasso Least Angle Regression	0.8003	1.1524	1.0729	0.0119	0.5008	0.18	0.012
ada	AdaBoost Regressor	0.8867	1.1618	1.0775	0.0027	0.4851	0.231	0.023
omp	Orthogonal Matching Pursuit	0.8072	1.1656	1.079	0.0006	0.5029	0.1812	0.01
dummy	Dummy Regressor	0.8107	1.1683	1.0803	-0.0016	0.5033	0.1827	0.01
knn	K Neighbors Regressor	0.7091	1.2439	1.1143	-0.0656	0.5012	0.1544	0.017
huber	Huber Regressor	0.7665	1.4728	1.2128	-0.2629	0.538	0.1499	0.039
dt	Decision Tree Regressor	0.6865	1.9758	1.405	-0.6979	0.6495	0.1604	0.016
par	Passive Aggressive Regressor	1.3615	4.9985	1.8137	-3.1395	0.6204	0.3573	0.013

Fig 3: Precision for Every Algorithm

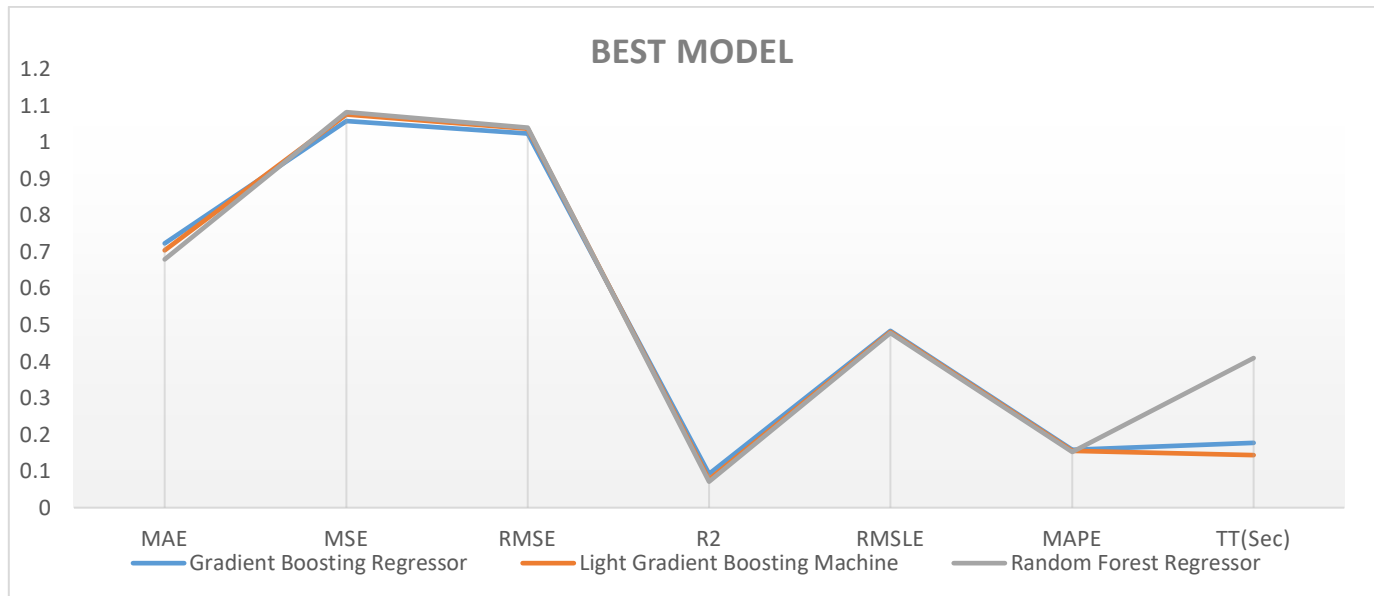


Fig 4: Best Model Analysis

5. CONCLUSION

In conclusion, the realm of data preprocessing in machine learning stands on the brink of transformative change. The challenges and inefficiencies identified, ranging from imbalanced class distributions to suboptimal hyperparameter optimization, not only present formidable obstacles but also offer compelling opportunities for innovation. This architectural analysis emphasizes the pivotal role of robust solutions in overcoming these hurdles and propelling the field of machine learning forward. The proposed AutoML system emerges as a pioneering and comprehensive approach to automating data preprocessing and model deployment. Its design aims to simplify the complexities of data preparation while upholding analytical rigor. Tailored components within the system address specific challenges such as imbalanced data management, hyperparameter optimization, feature engineering, and time-series handling. By offering a unified and integrated solution, this AutoML system has the potential to revolutionize data preprocessing in machine learning, providing practitioners with a powerful tool to enhance the quality and efficiency of their modeling endeavors.

FUTURE WORK:

In future developments, addressing the intricacies of data preprocessing is committed to several key areas of focus. Firstly, there is a commitment to enhancing automation capabilities, aiming to make the entire machine learning pipeline more accessible and efficient. This involves further streamlining processes, reducing manual intervention, and improving the overall user experience. Additionally, a focus on advanced model explainability is prioritized, aiming to enhance or refine methods that allow users to gain deeper insights into the decisions and predictions made by machine learning models. This is crucial for fostering trust and understanding in the application of these models.

REFERENCES

- [1] K. Goyle, Q. Xie, & V. Goyle, "DataAssist: A Machine Learning Approach to Data Cleaning and Preparation," eprint arXiv:2307.07119, 2023.
- [2] S. Juddoo, "Investigating Data Repair steps for EHR Big Data," in International Conference on Next Generation Computing Applications, 2022.

-
- [3] P. Ribeiro, P. Orzechowski, J. B. Wagenaar, & J. H. Moore, "Benchmarking AutoML algorithms on a collection of synthetic classification problems," eprint arXiv:2212.02704, 2022.
- [4] M. Abdelaal, C. Hammacher, & H. Schoening, "REIN: A Comprehensive Benchmark Framework for Data Cleaning Methods in ML Pipelines," eprint arXiv:2302.04702, 2023.
- [5] F. Neutatz, B. Chen, Y. Alkhatib, J. Ye, & Z. Abedjan, "Data Cleaning and AutoML: Would an Optimizer Choose to Clean?" Eprint Springer s13222-022-00413-2, 2022.
- [6] M. Abdelaal, R. Koparde, & H. Schoening, "AutoCure: Automated Tabular Data Curation Technique for ML Pipelines," eprint arXiv:2304.13636, 2023.
- [7] S. Holzer & K. Stockinger, "Detecting errors in databases with bidirectional recurrent neural networks," OpenProceedings ZHAW, 2022.
- [8] P. Li, Z. Chen, X. Chu, & K. Rong, "DiffPrep: Differentiable Data Preprocessing Pipeline Search for Learning over Tabular Data," eprint arXiv:2308.10915, 2023.
- [9] M. Singh, J. Cambronoero, S. Gulwani, V. Le, C. Negreanu, & G. Verbruggen, "DataVinci: Learning Syntactic and Semantic String Repairs," eprint arXiv:2308.10922, 2023.
- [10] S. Guha, F. A. Khan, J. Stoyanovich, & S. Schelter, "Automated Data Cleaning Can Hurt Fairness in Machine Learning-based Decision Making," in IEEE 39th International Conference on Data Engineering, 2023.
- [11] R. Wang, Y. Li, & J. Wang, "Sudowoodo: Contrastive Self-supervised Learning for Multi-purpose Data Integration and Preparation," eprint arXiv:2207.04122, 2022.
- [12] B. Hilprecht, C. Hammacher, E. Reis, M. Abdelaal, & C. Binnig, "DiffML: End-to-end Differentiable ML Pipelines," eprint arXiv:2207.01269, 2022.
- [13] V. Restat, M. Klettke, & U. Störl, "Towards a Holistic Data Preparation Tool," in EDBT/ICDT Workshops, 2022.
- [14] M. Nashaat, A. Ghosh, J. Miller, & S. Quader, "TabReformer: Unsupervised Representation Learning for Erroneous Data Detection," eprint <https://doi.org/10.1145/3447541>, 2021.
- [15] F. Calefato, L. Quaranta, F. Lanubile, & M. Kalinowski, "Assessing the Use of AutoML for Data-Driven Software Engineering," eprint arXiv:2307.10774, 2023.