# PIFE: PROGRESSIVE INSIGHT DRIVEN FEATURE ENGINEERING VIA MULTIMODAL REASONING

**Anonymous authors**Paper under double-blind review

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032033034

035

037

040

041

042

043

044

046 047

048

051

052

#### **ABSTRACT**

Despite significant advances in Automated Machine Learning (AutoML), one of its persistent blind spots remains the automation of data-centric tasks such as exploratory data analysis (EDA), contextual insight extraction, and feature engineering. These steps-often more critical than model selection itself-are still largely manual, domain-specific, and reliant on human intuition. Existing automated feature engineering (AutoFE) techniques either rely on rigid transformation sets or complex optimization strategies that struggle with interpretability and fail to leverage the rich, visual cues that guide human decision-making. In this work, we introduce PIFE: Progressive Insight driven Feature Engineering via Multimodal Reasoning; a novel AutoFE framework that employs multimodal language models as collaborative agents in an iterative pipeline. PIFE systematically performs automated EDA, generating statistical summaries and visualizations that are jointly interpreted through text-vision reasoning. These multimodal insights inform the synthesis of candidate transformations, represented as symbolic programs in executable Python code to ensure interpretability and reproducibility. By coupling iterative insight extraction with validation-driven refinement, PIFE produces highquality, interpretable features that consistently enhance the performance of diverse predictive models, outperforming existing AutoFE baselines. Extensive experiments across diverse tabular datasets demonstrate the effectiveness and adaptability of our approach, paving the way for a new class of human-aligned, insightaware AutoFE systems.

#### 1 Introduction

The rapid evolution of automated machine learning (AutoML) has significantly advanced model selection, hyperparameter tuning, and performance optimization (Chopde et al., 2025; Aragão et al., 2025; Hutter et al., 2019; Feurer et al., 2015; Olson & Moore, 2016; Erickson et al., 2020). However, AutoML tools continue to face limitations in automating data engineering tasks, particularly exploratory data analysis (EDA), feature insight extraction, and systematic feature engineering. These data-centric activities often dominate real-world machine learning workflows, where the transformation of raw tabular data into meaningful representations is a bigger bottleneck than model fitting. Although automated feature engineering (AutoFE) has emerged as a subfield within AutoML, traditional methods, such as expansion-reduction algorithms (Kanter & Veeramachaneni, 2015; Lam et al., 2021; Kaul et al., 2017; Shi et al., 2020; Katz et al., 2016) typically construct large search spaces composed of manually defined transformation operations and employ various search or optimization strategies to identify effective features. However, these methods are often limited by the rigidity of their predefined operations and generally lack the integration of domain-specific knowledge (Zhang et al., 2023).

To reduce the cost of searching through large feature space and generate data-driven features, learning-based AutoFE methods are proposed (Khurana et al., 2018; Nargesian et al., 2017; Chen et al., 2019; Zhu et al., 2022). However, these methods fall short in incorporating domain expertise and contextual insights from data exploration. Similarly, evolutionary methods focus on optimization strategies but neglect the nuanced, often visual cues that inform human-driven feature creation. Language-powered systems like CAAFE (Hollmann et al., 2023) and LLM-FE (Abhyankar et al., 2025) have shown promise in bridging this gap by generating candidate features based on dataset context and iterative refinement. However, these methods remain limited by feature simplicity, a

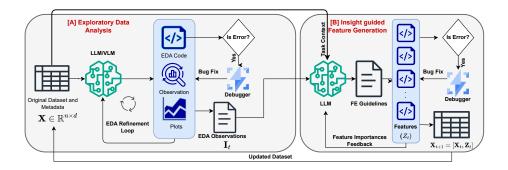


Figure 1: Overview of PIFE Framework. For a given dataset, PIFE goes through the following steps. (a) Exploratory Data Analysis and Data Insight Generation(b) Feature Generation via Symbolic Program Synthesis

lack of interpretability regarding why certain features should be created (as opposed to merely explaining what they represent), and the absence of a truly data-driven approach. Furthermore, visual patterns-such as distributional anomalies, multivariate correlations, or interaction structures-remain underutilized despite their centrality in manual feature engineering workflows.

This gap highlights the opportunity to harness recent advancements in Large Language Models that understand both textual and visual modalities to build a framework for automated insight extraction and feature engineering. These models are capable of interpreting not only data descriptions but also visualizations such as histograms, scatter plots, and heatmaps; elements that humans frequently rely on during feature engineering in real-world scenarios. Yet, their potential in automating feature generation grounded in rich exploratory insights has remained largely unexplored. A more effective AutoFE pipeline must seamlessly incorporate insights from both narrative and visual representations of data.

To address these challenges, we propose a novel AutoFE framework that integrates iterative EDA cycles using a unified reasoning engine capable of understanding both text and plots. Our system performs repeated rounds of insight extraction to build a deeper and comprehensive contextual understanding of the dataset, which then guides feature generation. The generated candidate features are evaluated using a downstream predictive model, where the corresponding feature importance serves as feedback to subsequent feature generation cycles. We argue that the broader process of feature engineering can be naturally decomposed into two complementary stages: (i) feature generation and (ii) feature selection. While the former aims to enrich the feature space, the latter plays a critical role in filtering redundant or irrelevant features and selecting an optimal subset that maximizes task performance. To emphasize the importance of this selection step, we conduct extensive experiments comparing diverse feature selection methodologies and demonstrate that incorporating effective selection strategies can further enhance the performance of automated feature engineering (AutoFE) pipelines. This feedback-driven, context-rich process enhances automation and interpretability, while aligning closely with the iterative and insight-informed nature of human data science workflows.

**Contributions.** The key contributions of this work are as follows:

- We propose the first automated feature engineering framework that integrates textual and visual exploratory data insights into a unified, iterative pipeline.
- We highlight the central role of feature selection in AutoFE by conducting extensive experiments across diverse selection methodologies, showing that effective selection strategies further boost both predictive performance and interpretability compared to state-of-the-art AutoFE methods.
- We conduct extensive experiments across various tabular datasets, demonstrating superior performance and enhanced interpretability compared to state-of-the-art AutoFE methods.

## 2 RELATED WORK

110 111 112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

108

109

Automated feature engineering (AutoFE) has emerged as a critical component in simplifying the model development pipeline by transforming raw data into informative representations. Early efforts such as Deep Feature Synthesis (DFS) (Kanter & Veeramachaneni, 2015), LFE (Nargesian et al., 2017), Cognito (Khurana et al., 2016), AutoFeat (Horn et al., 2020), and OpenFE (Zhang et al., 2023) employed exhaustive enumeration or heuristic-based transformation strategies, often relying on predefined operator sets and lacking semantic understanding of domain-specific relationships. OpenFE extended traditional methods via an expansion-reduction framework with incremental feature boosting and pruning, achieving strong empirical performance but still limited by its lack of contextual reasoning and domain adaptivity.

More recent methods address these limitations through learning-based strategies. TransGraph (Khurana et al., 2018), Neural Feature Search (NFS) (Chen et al., 2019), and DIFER (Zhu et al., 2022) adopted reinforcement learning and differentiable architecture search to explore high-dimensional transformation spaces more efficiently. DIFER, in particular, proposed a differentiable encoder-predictor-decoder pipeline to optimize feature embeddings in continuous space, though it primarily supports numerical transformations. FETCH (Li et al., 2023) approached AutoFE as a Markov Decision Process, using a policy network trained across datasets to learn transferable feature construction policies. Despite its generalizability, FETCH suffers from sparse rewards and computational overhead, echoing challenges seen in DIFER and NFS.

In parallel, large language models (LLMs) have shown promise in data-centric applications, leveraging their contextual understanding to perform data wrangling, imputation, and semantic reasoning over tabular data (Hegselmann et al., 2023; Narayan et al., 2022; Vos et al., 2022). CAAFE (Hollmann et al., 2023) was among the first to explore LLM-driven feature engineering, generating features based on dataset metadata and producing humanreadable descriptions. However, it lacks iterative feedback from prior search histories and relies on column descriptions to create features. OCTree (Nam et al., 2024) augmented this by incorporating decisiontree reasoning into LLM prompts, offering structured, contextual feedback for subsequent feature generation. OCTree performs iterative refinement of feature generation rules until improvements in downstream performance. While effective, this approach is susceptible to poor initialization in LLMs, which can hinder convergence and overall effectiveness.

LLM-FE (Abhyankar et al., 2025) takes a different approach by casting feature engineering as a program synthesis problem.

**Algorithm 1** PIFE: Insight Driven Iterative Feature Generation

**Input:** Dataset  $\mathbf{X}_0 \in \mathbb{R}^{n \times d}$ , target  $\mathbf{y}$ , context  $\mathcal{P}$  **Parameters:** Max EDA rounds K, Max iterations  $T_f$ , feature cap N

```
Output: Candidate feature sets \{\mathbf{Z}_t\}_{t=1}^{T_f}
   1: t \leftarrow 0, \mathbf{X}_0 \leftarrow \mathbf{X}
   2: while t < T_f do
                 while t < T_f do \mathcal{D}_t = \{(c_i, \tau_i)\}_{i=1}^{d_t} \quad \{\text{col. names } c_i, \text{ types } \tau_i\} \mathcal{C}_t^{(0)} = \mathcal{D}_t \cup \mathcal{P} for k = 1 to K do S_t^{(k)} = \text{Stats}(\mathbf{X}_t) \quad \{\text{skew, quantiles, corr.}\} \mathbf{I}_t^{(k)} = \text{LLM}(\mathcal{C}_t^{(k-1)} \cup S_t^{(k)}) \quad \{\text{stat. insights}\} \mathcal{V}_t^{(k)} = \text{Viz}(\mathbf{X}_t) \quad \{\text{plots}\} \mathbf{I}_{vlm,t}^{(k)} = \text{VLM}(\mathcal{V}_t^{(k)}) \quad \{\text{visual insights}\} \mathcal{C}_t^{(k)} = \mathcal{C}_t^{(k-1)} \cup \mathbf{I}_t^{(k)} \cup \mathbf{I}_{vlm,t}^{(k)}
                   end for \mathbf{I}_t = \bigcup_{k=1}^K (\mathbf{I}_t^{(k)} \cup \mathbf{I}_{vlm,t}^{(k)})
11:
12:
                     \mathcal{G}_t = \text{LLM}(\mathbf{I}_t)
13:
                                                                                                  {transformation rules}
                    \mathbf{Z}_t = \mathrm{LLM}(\mathcal{G}_t)
14:
                                                                                                            {feature programs}
15:
                    \hat{y}_t = \mathcal{M}(\mathbf{X}_t \cup \mathbf{Z}_t) {downstream evaluation}
                    \phi_t = \operatorname{Importance}(\mathcal{M}, \mathbf{Z}_t)
16:
                                                                                                                                    {feat. imp.}
                    \mathcal{P} \leftarrow \mathcal{P} \cup \phi_t
17:
                                                                                     {update feedback context}
                    t \leftarrow t + 1
18:
```

It combines LLM reasoning, evolutionary strategies, and memory buffers to maintain a population of candidate features, using both validation scores and information-theoretic feedback to guide selection. This method addresses sparse rewards and brittleness in prompting; it introduces a new challenge: by conditioning future generations on previously successful feature transformations, the model may become biased toward certain transformation patterns. This can skew the process toward exploitation, limiting its ability to explore novel and potentially better features.

19: **end while=**0

Our framework leverages LLMs as agents for in-depth EDA, identifying outliers, feature interactions, and distributional patterns, which are mapped to candidate transformations and validated

via downstream performance-balancing automation, interpretability, and domain adaptivity beyond prior AutoFE methods.

## 3 METHODOLOGY

Our framework leverages Large Language Models (LLMs) and Vision-Language Models (VLMs) as autonomous reasoning agents in an iterative feature engineering process. By combining semantic understanding, language-guided program synthesis, and feedback from downstream predictive models, the system emulates a data scientist's workflow to generate high-quality, interpretable features. Section 3.1 formalizes the problem, Section 3.1.2 describes the iterative insight-driven feature generation, and Section 3.2 details the integration of downstream feedback across iterations.

#### 3.1 PROBLEM FORMULATION

Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be a tabular dataset with n samples and d initial features, and  $\mathbf{y} \in \mathbb{R}^n$  the corresponding target vector. Our objective is to construct an enriched feature set  $\mathbf{X}^* \in \mathbb{R}^{n \times d^*}$ , where  $d^* \geq d$ , by iteratively generating candidate feature transformations.

At iteration  $t \in \{1, ..., T_f\}$ , the framework produces candidate features  $\mathbf{Z}_t$  guided by insights  $\mathbf{I}_t$ , evaluates them with a predictive model  $\mathcal{M}$ , and updates a feedback context  $\mathcal{P}$  containing feature importances:

$$\hat{\mathbf{y}}_t = \mathcal{M}(\mathbf{X}_t \cup \mathbf{Z}_t), \quad \phi_t = \text{Importance}(\mathcal{M}, \mathbf{Z}_t), \quad \mathcal{P} \leftarrow \mathcal{P} \cup \phi_t$$

The feature space is then incrementally updated:  $\mathbf{X}_{t+1} = [\mathbf{X}_t, \mathbf{Z}_t]$ . After  $T_f$  iterations, the final feature matrix  $\mathbf{X}^*$  is selected as the one maximizing cross-validation performance:

$$\mathbf{X}^* = \arg \max_{\mathbf{X}_t} \text{CV\_Score}(\mathcal{M}, \mathbf{X}_t, \mathbf{y})$$
 (1)

We mimic a data scientist's flow to perform Exploratory Data Analysis (EDA) on the data by generating various sorts of analysis, such as feature distribution analysis, interaction analysis - feature pairs and non-linear relationships, temporal categorical analysis - temporal trends and categorical encodings, etc. which help explore complex data patterns and relationships among data points as shown in Figure 2. These insights, grounded in the underlying datasets, contain meaningful information for feature generation. Paired with LLMs as optimizers, these insights guide the way for feature transformation rules. Feature generation proceeds in two stages: (1) iterative extraction of insights using LLMs and VLMs, and (2) symbolic feature program synthesis guided by these insights.

#### 3.1.1 EDA INSIGHT EXTRACTION VIA LLM/VLM

At each iteration t, the system constructs a structured dataset description  $\mathcal{D}_t = \{(c_i, \tau_i)\}_{i=1}^{d_t}$ , where  $c_i$  is a column name and  $\tau_i$  is its type, combined with task metadata and domain context  $\mathcal{P}$ :

$$\mathcal{C}_t^{(0)} = \mathcal{D}_t \cup \mathcal{P}$$

The LLM performs iterative reasoning across K internal rounds; at each round k, it computes additional statistics  $S_t^{(k)}$  (e.g., quantiles, skewness, missingness), generates updated insights  $\mathbf{I}_t^{(k)} = \text{LLM}(\mathcal{C}_t^{(k-1)} \cup S_t^{(k)})$ , and expands the context  $\mathcal{C}_t^{(k)} = \mathcal{C}_t^{(k-1)} \cup \mathbf{I}_t^{(k)}$ .

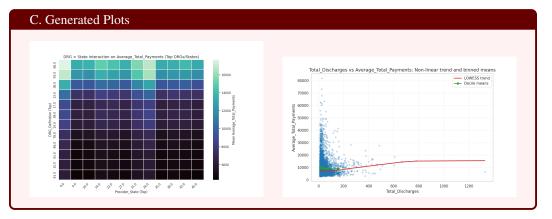
In parallel, visual summaries  $\mathcal{V}_t^{(k)}$  (plots, distributions) are analyzed with a VLM, producing visual insights  $\mathbf{I}_{\text{vlm},t}^{(k)} = \text{VLM}(\mathcal{V}_t^{(k)})$ . Consolidated insights for iteration t are:

$$\mathbf{I}_t = igcup_{k=1}^K (\mathbf{I}_t^{(k)} \cup \mathbf{I}_{\mathrm{vlm},t}^{(k)})$$

## A. Rationale for EDA

To study interaction and non-linear effects, I examine the joint surface of Average\_Covered\_Charges and reimbursement\_rate against the target to detect curvature and interaction; cross-category heterogeneity via a DRG-by-State heatmap to expose geographic—clinical interplay; and utilization effects by relating Total\_Discharges to the target with a smooth trend plus binned means to reveal non-linear scaling and heteroskedasticity. Expected outputs highlight regimes where combinations of drivers jointly shift payment levels, and where certain DRG—State combinations systematically differ. Limitation: heavy tails may distort binning; quantile-based bins and downsampling are applied as needed.

## B. EDA Code Snippet



#### D. EDA Analysis

**Total Discharges vs target scatter with LOWESS:** The relationship is non-linear with a concave (increasing-then-flattening) trend; variance is high at low volumes and shrinks as discharges grow. This motivates monotone, variance-stabilizing transforms and regime features (e.g., ranks/bins).

**DRG** × **State heatmap:** There is strong cross-category heterogeneity—within the same DRG, Provider\_State causes sizable shifts in mean payments, and DRGs also differ markedly in their baseline level. This supports both main effects (DRG, State) and their interaction.

Figure 2: Exemplery run of PIFE on medical\_charges\_nominal dataset showing the process of generating data insights. First, the rationale is generated, creating a plan for exploratory analysis to be conducted. In B, this plan is translated into a program for EDA. In C, when the code is executed, analysis plots are generated, and at the end, plots and statistics are analyzed to generate statistics.

## 3.1.2 LLM-GUIDED FEATURE PROGRAM GENERATION

Given the insight context  $I_t$ , the LLM produces a set of transformation guidelines  $\mathcal{G}_t = \{g_{t,1}, \dots, g_{t,n_t}\}$ . Each guideline is converted into a symbolic feature program and Reverse Pol-

ish Notation (RPN) representation, yielding generated features:

$$\mathbf{Z}_{t} = \{z_{t,1}, \dots, z_{t,n_t}\}, \quad n_t \le N, \quad z_{t,i} = \text{LLM}(g_{t,i})$$

These features are appended to the current dataset,  $\mathbf{X}_t \leftarrow [\mathbf{X}_{t-1}, \mathbf{Z}_t]$ .

#### 3.2 ITERATIVE REFINEMENT AND EVALUATION

We select all the generated features, and the respective feature importances computed after each iteration provide context  $\mathcal{P}$  for subsequent feature generation as well as insight extraction. This iterative loop ensures that newly generated features complement the existing space, help avoid duplicate features, and consistently improve predictive performance.

$$\mathbf{X}_{t+1} = [\mathbf{X}_t, \mathbf{Z}_t], \quad \phi_t = \text{Importance}(\mathcal{M}, \mathbf{Z}_t), \quad \mathcal{P} \leftarrow \mathcal{P} \cup \phi_t$$

After  $T_f$  iterations, the final feature set  $\mathbf{X}^*$  is chosen based on cross-validation performance of the downstream model as mentioned in equation 1.

#### 4 EXPERIMENTS

In this section, we evaluate PIFE over several classification and regression datasets spanning across various domains such as healthcare, finance, real estate, weather forecasting, etc. Our experiments reveal that PIFE consistently improves the performance of predictive models (Section 4.2). Ablation studies (Section 4.3) show that data-grounded insight extraction helps create features that are more aligned to the downstream objective. Also, feature selection is often not focused on in the scope of feature engineering, which plays a pivotal role in boosting the performance of predictive models.

#### 4.1 EXPERIMENTAL SETUP

**Datasets.** We evaluate PIFE across 22 tabular tasks, encompassing both classification and regression objectives. The majority of datasets are drawn from prior AutoFE literature (Li et al., 2023), ensuring coverage of diverse domains, scales, and complexity levels. Additionally, we include a set of recent Kaggle datasets (Kaggle) (e.g., ps5\_episode\_3, ps5\_episode\_4), which were not part of the pretraining corpora of large language models, providing a test of robustness to novel data sources. Detailed dataset information is provided in Table 6.

**Metrics.** For classification tasks, we use F1-micro (Sokolova & Lapalme, 2009), and for regression tasks, we use (1 - relative absolute error) (Shcherbakov et al., 2013) as the evaluation metric for downstream models. Higher values correspond to better model performance. To quantify improvements, we also report the percentage increase over a baseline score, reflecting relative efficacy.

Baselines. We compare PIFE against a diverse set of baseline methods representing key paradigms in automated feature engineering, all of which have publicly available, executable open-source implementations to ensure reproducibility. Heuristic-based approaches include AutoFeat Horn et al. (2020), DFS (Deep Feature Synthesis) Kanter & Veeramachaneni (2015), and OpenFE Zhang et al. (2023), which rely on expansion and reduction strategies over predefined transformations. Among LLM-based approaches, we include CAAFE Hollmann et al. (2023), leveraging LLMs for feature generation and refinement using metadata, prompts, or reasoning frameworks, and OCTree Nam et al. (2024), which employs rule-based feature generation and CART decision tree inputs to improve feature quality. The CAAFE and OCTree implementations were adapted to support newer LLM models and extended to handle both classification and regression tasks, with additional metrics introduced for fairer comparison. However, certain recent methods, such as LLM-FE Abhyankar et al. (2025), are excluded due to incomplete publication of methodology and evaluation details, which would limit fair comparison. Additional discussion is provided in Appendix A.6.

**Implementation Details.** To ensure reliable evaluation, we perform 5-fold cross-validation on the training set, mitigating overfitting and yielding robust performance estimates. Results are reported as mean  $\pm$  standard deviation over three random seeds (42, 44, 46) to account for stochasticity in LLMs and training pipelines. We use the given specific versions of LLMs: gpt-4.1-2025-04-14 and gpt-5-2025-08-07. For fairness, all datasets are preprocessed by imputing or removing missing values and encoding categorical variables, as most downstream models lack native support. Further details on LLM prompting strategies (Appendix A.8), hyperparameters (Appendix A.5), and additional configuration settings are provided in Appendix A.

Table 1: Comparison of AutoFE methods across method compatibility and performance (mean  $\pm$  std) for different LLMs. Tick ( $\checkmark$ ) indicates presence, cross ( $\checkmark$ ) indicates absence of a feature. Results are averaged across 3 seeds, with each seed evaluated using 5-fold cross-validation and Random Forest as the predictive model. We report the f1-micro score for classification and (1-rae) for regression datasets.

Method	Context Aware	Without Description	Interpretable Feature	LLM	Avg. Score (%)
Baseline	X	×	X	-	$0.7558 \pm 0.1017$
DFS	X	✓	X	-	$0.7718 \pm 0.1101  (2.12\%)$
Autofeat	X	✓	Х	-	$0.7651 \pm 0.0948  (1.23\%)$
OpenFE	X	✓	X	-	$0.7684 \pm 0.0922  (1.67\%)$
CAAFE	<b>✓</b>	×	✓	gpt-4.1* gpt-5*	$0.7791 \pm 0.1068 (3.08\%)$ $0.7900 \pm 0.0996 (4.53\%)$
OCTree	X	<b>✓</b>	×	gpt-4.1* gpt-5*	$0.7745 \pm 0.0958 $ (2.47%) $0.7750 \pm 0.0969 $ (2.54%)
PIFE (Ours)	<b>✓</b>	✓	✓	gpt-4.1* gpt-5*	$0.7908 \pm 0.1105  (4.63\%) \ 0.7917 \pm 0.1037  (4.75\%)$

<sup>\*</sup> All experiments were conducted using LLM versions gpt-4.1-2025-04-14 and gpt-5-2025-08-07.

Table 2: Performance comparison of Baseline, w/o EDA, and w/ EDA across multiple datasets under the No Feature Selection setting. The best value per dataset is highlighted in bold. Percentage improvement over the baseline is shown in parentheses. We report f1-micro score for classification (\*) and (1-relative absolute error) for regression (†) datasets.

Dataset	Baseline	w/o EDA	w/ EDA
adult*	0.850	0.852 (0.3%)	0.853 (0.3%)
fertility*	0.829	0.873 (5.4%)	0.880 (6.2%)
medical_charges_nominal <sup>†</sup>	0.891	0.922 (3.5%)	0.903 (1.4%)
openml_586 <sup>†</sup>	0.613	0.729 (18.8%)	0.772 (25.9%)
pima_indian*	0.700	0.759 (8.3%)	0.754 (7.7%)
ps5_episode_4 <sup>†</sup>	0.571	0.577 (1.2%)	0.578 (1.3%)

#### 4.2 Performance Comparisons

Table 1 highlights PIFE as the most effective and practical AutoFE method: it attains the top average score while preserving semantic interpretability and context awareness, and it works even without dataset descriptions. Classical baselines (DFS, Autofeat, OpenFE) offer modest gains but lack contextual understanding and interpretability. Among LLM-based methods, PIFE leads fairly: with gpt-4.1, it improves 4.63% over the Baseline (without feature engineering) versus 3.08% for CAAFE; with gpt-5, the gap narrows, but PIFE still edges ahead, suggesting stronger reasoning models reduce, but do not erase, method-level differences. These results demonstrate that PiFE consistently improves performance across seeds and folds, producing interpretable, context-aware features with minimal dependence on the latest LLMs. Full results are reported in Appendix Table 7.

#### 4.3 ABLATIONS

Impact of EDA. To assess the contribution of the EDA component in insight-driven feature generation, we compare model performance with and without EDA in Table 2. Even without EDA, the generated features are optimized and achieve competitive results. However, EDA provides a data-grounded mechanism for feature generation, enabling the capture of complex relationships and trends that are difficult to model when relying solely on LLM optimization or metric-based feedback. We observe that certain datasets exhibit limited or no performance gains from the inclusion of EDA. This may occur when the original features already capture sufficient signal, the dataset is small, or added features introduce redundancy rather than value.

Feature Selection: Trade-offs and Downstream Performance. Although PIFE can generate interpretable and statistically strong features, they might not reflect the same way on performance. Therefore, effective feature selection is crucial to identify the optimal subset of features. To evaluate this, we compare three approaches: Model-based Feature Importance(MFI), Conditional Mutual Information-based Bayesian Optimization (CMI-BO), and Genetic Algorithm(GA) (see the Appendix A.7). CMI-BO and MBFI did not show any improvement in performance. The genetic algorithm, while computationally intensive, gave the best performance compared to others. We can see that the set of features generated by PIFE enabled good exploration inthe Genetic Algorithm. In this way, PIFE complements these approaches by producing a concise, high-quality subset of features, making downstream optimization more efficient while maintaining strong predictive performance.

Table 3: Performance comparison of PIFE with feature selection during run (using CMI and validation) and after the run (using a Genetic Algorithm), with GPT-5 as the downstream model. Reported metrics are F1-micro for classification datasets and (1-relative absolute error) for regression datasets.

Dataset	Baseline	PIFE	PIFE(CMI-BO)	PIFE(MFI)	PIFE(GA)
adult	$0.851 \pm 0.001$	$0.853 \pm 0.001$	$0.853 \pm 0.001$	$0.851 \pm 0.002$	$0.855 \pm 0.001$
fertility	$0.853 {\pm} 0.006$	$0.88 {\pm} 0.010$	$0.873 \pm 0.015$	$0.873 \pm 0.006$	$0.89 {\pm 0.008}$
openml_586	$0.662 \pm 0.009$	$0.765 {\pm} 0.012$	$0.742 \pm 0.025$	$0.760 \pm 0.001$	$\textbf{0.792} \!\pm\! 0.002$
pima_indian	$0.739 {\pm} 0.010$	$0.754 \pm 0.01$	$0.755 \pm 0.005$	$0.76 \pm 0.010$	$0.777 \!\pm\! 0.002$

Table 4: Feature Transferability of PiFE-generated Features to Deep Learning Models (MLP, TabPFN(Hollmann et al., 2023), and HyperFast (Bonet et al., 2024)). \* denotes classification and †denotes regression tasks. HyperFast (NA) only runs on classification tasks.

Dataset	M	LP	Tab	PFN	HyperFast	
	Baseline	PIFE (Ours)	Baseline	PIFE (Ours)	Baseline	PIFE (Ours)
hepatitis*	0.862±0.020	0.852±0.016	0.832±0.005	0.826±0.005	0.815±0.015	0.843±0.007
airfoil <sup>†</sup>	$0.735 \pm 0.001$	$0.802 {\pm} 0.001$	$0.955 {\pm} 0.000$	$0.952 \pm 0.001$	NA	NA
credit_approval*	$0.888 {\pm} 0.002$	$0.884 {\pm} 0.002$	$0.874 \pm 0.005$	$0.863 \pm 0.005$	$0.861 \pm 0.006$	$0.847 \pm 0.003$
spectf*	$0.828 {\pm} 0.003$	$0.81 \pm 0.014$	$0.798 \pm 0.005$	$0.806 {\pm} 0.005$	$0.792 \pm 0.007$	$0.805 {\pm 0.01}$
megawatt_1*	$0.900 \pm 0.008$	$0.908 \pm 0.004$	$0.893 \pm 0.003$	$0.895 {\pm} 0.008$	$0.865 {\pm} 0.008$	$0.884 {\pm} 0.007$
housing_boston <sup>†</sup>	$0.701 \!\pm\! 0.001$	$\textbf{0.706} \!\pm\! \textbf{0.001}$	$0.873 \!\pm\! 0.002$	$\textbf{0.875} {\pm 0.000}$	NA	NA

Feature transferability is critically dependent on model inductive bias. Transformations that encode tree-like, thresholding behaviour generally transfer well to tree ensembles but can degrade performance for neural or transformer architectures that favour smooth, continuous representations or learned embeddings. Moreover, nearceiling baseline performance leaves little headroom for gains.

Table 5: Performance comparison of PIFE and PIFE $^{\dagger}$  (extended to OpenFE) across competitions. Values: mean  $\pm$  standard deviation. All results are based on gpt-5. f1-micro score for classification and (1-relative absolute error) for regression datasets.

Competition	PIFE	PIFE <sup>†</sup>
adult	$0.851 \pm 0.002$	$0.855 \pm 0.002$
fertility	$0.870 \pm 0.040$	$0.870 \pm 0.040$
medical_charges_nominal	$0.905 \pm 0.000$	$\textbf{0.907} \pm \textbf{0.001}$
openml_586	$0.773 \pm 0.023$	$\textbf{0.790} \pm \textbf{0.010}$
openml_607	$0.732 \pm 0.012$	$\textbf{0.752} \pm \textbf{0.023}$
ps5_episode_4	$0.578 \pm 0.001$	$\textbf{0.579} \pm \textbf{0.016}$
Average	$0.778 \pm 0.107$	$0.780 \pm 0.105$

Finally, transfer success is dataset-dependent, being modulated by sample size, noise, feature types, and the specific nature of engineered transformations. Engineered features should be validated on the intended downstream model family.

**Integrating with other AutoFE Methods.** Engineered features from PIFE can serve as input to other AutoFE frameworks. We experimented with OpenFE as the integrated framework and report the results in Table 5. Overall, integrating PIFE features with OpenFE didn't result in a

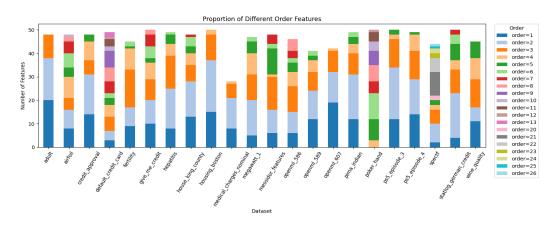


Figure 4: Order of features created per competition. This is based on the gpt-5 runs from Table 1

substantial improvement. This can be seen as a positive outcome: PIFE already identifies a strong set of features on its own. By leveraging insights from exploratory data analysis (EDA) and domain knowledge encoded in LLMs, along with natural language descriptions of the data, PIFE generates features that are both meaningful and predictive. Even after exploring a large space of additional candidate features ( $\sim$ 2000) in OpenFE, there is little to no gain, and in some cases, performance slightly decreases. This highlights the robustness and quality of the features generated by PIFE.

**Feature Order Analysis.** PIFE is capable of creating higher-order meaningful features and can be scaled well for creating complex interactions between features due to its iterative nature of feature generation. We can see that some of the new features created have an order greater than the number of feature engineering steps. From this, we can infer that LLM is attempting to create complex, higher-order features in a single feature engineering step.

**Flexibility Towards Predictive Models.** PiFE generalizes well across various downstream predictive models, showing consistent improvements with Logistic Regression, Random Forest, and XGBoost. As expected, tree-based models, which can capture non-linear relationships among features, generally outperform linear models such as Logistic Regression.

#### 5 Conclusion

We present PIFE, a multimodal AutoFE framework that leverages textual and visual insights from datasets to iteratively generate and select predictive features. By combining exploratory data analysis with LLMguided reasoning, PIFE automates feature engineering in a way that mirrors human workflows, enhancing both performance and interpretability. Our experiments across diverse tabular datasets demonstrate that effective feature selection amplifies the benefits of automated feature generation. However, PIFE has limitations: large feature sets can

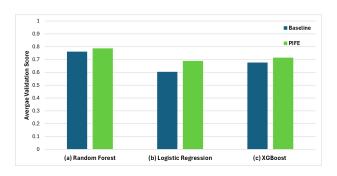


Figure 3: Performance of PIFE with predictive models. All results use GPT-5 as base LLM with the same hyperparameters as the main table. We report f1-micro for classification and (1-relative absolute error) for regression datasets.

lead to prompt size constraints, automatically derived dataset descriptions and visualizations may be noisy, and LLMs can produce plausible but ungrounded features. Additionally, results can vary across random seeds, underscoring the importance of multiple runs for robustness.

Future work includes improving multimodal reasoning, fine-tuning models for better feature generation, integrating human-in-the-loop interactions, and incorporating continuously updated

datasets to reduce memorization biases. PIFE represents a step toward context-aware, interpretable, and robust automated feature engineering.

#### REFERENCES

- Nikhil Abhyankar, Parshin Shojaee, and Chandan K Reddy. Llm-fe: Automated feature engineering for tabular data with llms as evolutionary optimizers. *arXiv preprint arXiv:2503.14434*, 2025.
- Marcelo VC Aragão, Augusto G Afonso, Rafaela C Ferraz, Rairon G Ferreira, Sávio G Leite, Felipe AP de Figueiredo, and Samuel B Mafra. A practical evaluation of automl tools for binary, multiclass, and multilabel classification. *Scientific Reports*, 15(1):17682, 2025.
- David Bonet, Daniel Mas Montserrat, Xavier Giró-i Nieto, and Alexander G Ioannidis. Hyperfast: Instant classification for tabular data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 11114–11123, 2024.
- Xiangning Chen, Qingwei Lin, Chuan Luo, Xudong Li, Hongyu Zhang, Yong Xu, Yingnong Dang, Kaixin Sui, Xu Zhang, Bo Qiao, et al. Neural feature search: A neural architecture for automated feature engineering. In 2019 IEEE International Conference on Data Mining (ICDM), pp. 71–80. IEEE, 2019.
- Abhishek Chopde, Fardeen Pettiwala, Sankar Kirubananth, Sai Kiran Botla, and Pachipulusu Ayyappa Kethan. PiML: Automated machine learning workflow optimization using LLM agents. In *AutoML 2025 Methods Track*, 2025. URL https://openreview.net/forum?id=NwlqBpsjZz.
- Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint arXiv:2003.06505*, 2020.
- Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. *Advances in neural information processing systems*, 28, 2015.
- Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David Sontag. Tabllm: Few-shot classification of tabular data with large language models, 2023. URL https://arxiv.org/abs/2210.10723.
- Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. *Advances in Neural Information Processing Systems*, 36:44753–44775, 2023.
- Franziska Horn, Robert Pack, and Michael Rieger. The autofeat python library for automated feature engineering and selection, 2020. URL https://arxiv.org/abs/1901.07329.
- Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods, systems, challenges.* Springer Nature, 2019.
- Kaggle. Kaggle: Your machine learning and data science community. https://www.kaggle.com. Accessed: 2025-09-24.
  - James Max Kanter and Kalyan Veeramachaneni. Deep feature synthesis: Towards automating data science endeavors. In 2015 IEEE international conference on data science and advanced analytics (DSAA), pp. 1–10. IEEE, 2015.
  - Gilad Katz, Eui Chul Richard Shin, and Dawn Song. Explorekit: Automatic feature generation and selection. In *2016 IEEE 16th international conference on data mining (ICDM)*, pp. 979–984. IEEE, 2016.
- Ambika Kaul, Saket Maheshwary, and Vikram Pudi. Autolearn—automated feature generation and selection. In 2017 IEEE International Conference on data mining (ICDM), pp. 217–226. IEEE, 2017.

- Udayan Khurana, Deepak Turaga, Horst Samulowitz, and Srinivasan Parthasrathy. Cognito: Automated feature engineering for supervised learning. In 2016 IEEE 16th international conference on data mining workshops (ICDMW), pp. 1304–1307. IEEE, 2016.
  - Udayan Khurana, Horst Samulowitz, and Deepak Turaga. Feature engineering for predictive modeling using reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
  - Hoang Thanh Lam, Beat Buesser, Hong Min, Tran Ngoc Minh, Martin Wistuba, Udayan Khurana, Gregory Bramble, Theodoros Salonidis, Dakuo Wang, and Horst Samulowitz. Automated data science for relational data. In 2021 IEEE 37th International Conference on Data Engineering (ICDE), pp. 2689–2692. IEEE, 2021.
  - Liyao Li, Haobo Wang, Liangyu Zha, Qingyi Huang, Sai Wu, Gang Chen, and Junbo Zhao. Learning a data-driven policy network for pre-training automated feature engineering. In *The Eleventh International Conference on Learning Representations*, 2023.
  - Octavio César Mesner and Cosma Rohilla Shalizi. Conditional mutual information estimation for mixed, discrete and continuous data. *IEEE Transactions on Information Theory*, 67(1):464–484, 2020.
  - Jaehyun Nam, Kyuyoung Kim, Seunghyuk Oh, Jihoon Tack, Jaehyung Kim, and Jinwoo Shin. Optimized feature generation for tabular data via llms with decision tree reasoning. Advances in Neural Information Processing Systems, 37:92352–92380, 2024.
  - Avanika Narayan, Ines Chami, Laurel Orr, Simran Arora, and Christopher Ré. Can foundation models wrangle your data?, 2022. URL https://arxiv.org/abs/2205.09911.
  - Fatemeh Nargesian, Horst Samulowitz, Udayan Khurana, Elias B Khalil, and Deepak S Turaga. Learning feature engineering for classification. In *Ijcai*, volume 17, pp. 2529–2535, 2017.
  - Randal S Olson and Jason H Moore. Tpot: A tree-based pipeline optimization tool for automating machine learning. In *Workshop on automatic machine learning*, pp. 66–74. PMLR, 2016.
  - Maxim Vladimirovich Shcherbakov, Adriaan Brebels, Nataliya Lvovna Shcherbakova, Anton Pavlovich Tyukov, Timur Alexandrovich Janovsky, Valeriy Anatol'evich Kamaev, et al. A survey of forecast error measures. *World applied sciences journal*, 24(24):171–176, 2013.
  - Qitao Shi, Ya-Lin Zhang, Longfei Li, Xinxing Yang, Meng Li, and Jun Zhou. Safe: Scalable automatic feature engineering framework for industrial tasks. In 2020 IEEE 36th International Conference on Data Engineering (ICDE), pp. 1645–1656. IEEE, 2020.
  - Marina Sokolova and Guy Lapalme. A systematic analysis of performance measures for classification tasks. *Information processing & management*, 45(4):427–437, 2009.
  - David Vos, Till Döhmen, and Sebastian Schelter. Towards parameter-efficient automation of data wrangling tasks with prefix-tuning. In *NeurIPS 2022 First Table Representation Workshop*, pp. 1–9, 2022.
  - Tianping Zhang, Zheyu Aqa Zhang, Zhiyuan Fan, Haoyan Luo, Fengyuan Liu, Qian Liu, Wei Cao, and Li Jian. Openfe: Automated feature generation with expert-level performance. In *International Conference on Machine Learning*, pp. 41880–41901. PMLR, 2023.
  - Guanghui Zhu, Zhuoer Xu, Chunfeng Yuan, and Yihua Huang. Difer: differentiable automated feature engineering. In *International Conference on Automated Machine Learning*, pp. 17–1. PMLR, 2022.
  - Yufeng Zou, Jean Utke, Diego Klabjan, and Han Liu. Automated feature engineering by prompting, 2025. URL https://openreview.net/forum?id=ZXO7iURZfW.

## A APPENDIX

#### A.1 REPRODUCIBILITY

We release our code to ensure reproducibility of experiments at https://anonymous.4open.science/r/pife-7FE6. The repository includes the main PiFE pipeline, integrations with baseline frameworks, and scripts for ablation studies. Additionally, we provide the datasets used in our experiments, including modifications made to publicly available datasets to facilitate consistent evaluation.

#### A.2 DATASET COLLECTION, PREPROCESSING AND RESULTS

The datasets are collected by ensuring that contamination is minimal. The formatted data description contains the task, dataset description, target, and objective. We are handling missing values by replacing them with numeric values using the mean and filling missing categorical values with the most frequent value. Categorical columns are encoded with ordinal encoding. The EDA process is informed about the preprocessing steps taken by providing the preprocessing context to the module. C and R in Table 6 represent Classification and Regression task types, respectively.

Table 6: Summary of Benchmark Datasets

Name	Alias	Source	Type (C/R)	Inst./Feat.	Subject Area
Adult	adult	UCI	С	48842 / 14	Social Science
Credit Approval	credit_approval	UCI	C	690 / 15	Business
Default Credit Card	default_credit_card	UCI	C	30000 / 23	Business
Fertility	fertility	UCI	C	100 / 9	Health and Medicine
Give Me Some Credit	give_me_credit	Kaggle	C	251503 / 10	Finance
Hepatitis	hepatitis	UCI	C	155 / 19	Health and Medicine
Megawatt1	megawatt_1	OpenML	C	253 / 37	Mathematics
Diabetic Retinopathy Debrecen	messidor_features	ÚČI	C	1151 / 19	Health and Medicine
PimaIndian	pima_indian	Kaggle	C	768 / 9	Health and Medicine
Poker Hand	poker_hand	UCI	C	1025010 / 10	Games
Rainfall Dataset	ps5_episode_3	Kaggle	C	2920 / 11	Weather and Climate
SPECTF	spectf	UCI	C	267 / 44	Health and Medicine
Statlog German Credit	statlog_german_credit	UCI	C	1000 / 20	Social Science
Wine Quality	wine_quality	UCI	C	4898 / 11	Business
Airfoil	airfoil	UCI	R	1503 / 5	Physics and Chemistry
House King County	house_king_county	Kaggle	R	21613 / 20	Business
Housing Boston	housing_boston	UCI	R	506 / 13	Finance
Medical Charges Nominal	medical_charges_nominal	OpenML	R	163065 / 11	Business
OpenML 586	openml_586	OpenML	R	1000 / 25	Mathematics
OpenML 589	openml_589	OpenML	R	1000 / 25	Mathematics
OpenML 607	openml_607	OpenML	R	1000 / 50	Mathematics
Podcast Listening Time	ps5_episode_4	Kaggle	R	1000000 / 10	Entertainment

Table 7: Performance comparison across different AutoFE methods. Values represent mean  $\pm$  standard deviation of the metric score. We use F1-micro for classification and (1 - rae) for regression tasks. Random Forest was used as a downstream predictive model.

Competition	Baseline	OpenFE	DFS	Autofeat	CA	AFE	OCT	REE	PIFE	(Ours)
	(NO FE)	I	l		gpt-4.1	gpt-5	gpt-4.1	gpt-5	gpt-4.1	gpt-5
adult	0.8511 ± 0.0007	0.8538 ± 0.0004	0.8547 ± 0.0009	0.8472 ± 0.0003	0.8532 ± 0.0021	0.8500 ± 0.0012	0.8547 ± 0.0008	0.8529 ± 0.0009	0.8853 ± 0.0609	0.8526 ± 0.0013
airfoil	0.7436 ± 0.0018	$0.7473 \pm 0.0097$	$0.7384 \pm 0.0037$	$0.7508 \pm 0.0042$	$0.7428 \pm 0.0013$	$0.7613 \pm 0.0024$	$0.7555 \pm 0.0062$	$0.7513 \pm 0.0045$	$0.7480 \pm 0.0047$	$0.7588 \pm 0.0060$
credit_approval	$0.8531 \pm 0.0093$	$0.8585 \pm 0.0051$	$0.8589 \pm 0.0068$	$0.8604 \pm 0.0080$	$0.8609 \pm 0.0088$	$0.8565 \pm 0.0124$	$0.8667 \pm 0.0025$	$0.8681 \pm 0.0063$	$0.8763 \pm 0.0344$	0.8659 ± 0.0010
default_credit_card	$0.8070 \pm 0.0007$	$0.8072 \pm 0.0007$	$0.8079 \pm 0.0008$	$0.8087 \pm 0.0021$	$0.8061 \pm 0.0015$	$0.8067 \pm 0.0013$	$0.8093 \pm 0.0009$	$0.8090 \pm 0.0002$	$0.8092 \pm 0.0008$	$0.8088 \pm 0.0003$
fertility	$0.8533 \pm 0.0058$	$0.8633 \pm 0.0115$	$0.8733 \pm 0.0058$	$0.8533 \pm 0.0058$	$0.8600 \pm 0.0100$	$0.8533 \pm 0.0058$	0.8900 ± 0.0000	$0.8900 \pm 0.0100$	$0.8733 \pm 0.0153$	$0.8800 \pm 0.0100$
give_me_credit	$0.9336 \pm 0.0001$	$0.9328 \pm 0.0002$	$0.9334 \pm 0.0004$	$0.9337 \pm 0.0002$	$0.9339 \pm 0.0005$	$0.9335 \pm 0.0002$	0.9341 ± 0.0001	$0.9342 \pm 0.0004$	$0.9337 \pm 0.0001$	$0.9337 \pm 0.0002$
hepatitis	$0.8151 \pm 0.0372$	$0.8172 \pm 0.0244$	0.7871 ± 0.0171	$0.8194 \pm 0.0171$	$0.8280 \pm 0.0037$	$0.8086 \pm 0.0074$	$0.8366 \pm 0.0099$	$0.8473 \pm 0.0207$	$0.8151 \pm 0.0134$	$0.8172 \pm 0.0325$
house_king_county	$0.6865 \pm 0.0015$	$0.6887 \pm 0.0016$	0.6792 ± 0.0011	0.6875 ± 0.0018	$0.6878 \pm 0.0037$	$0.6948 \pm 0.0034$	$0.6883 \pm 0.0008$	0.6897 ± 0.0013	0.6958 ± 0.0010	$0.6909 \pm 0.0014$
housing_boston	$0.6388 \pm 0.0063$	$0.6488 \pm 0.0108$	$0.6306 \pm 0.0024$	$0.6482 \pm 0.0033$	$0.6380 \pm 0.0100$	$0.6459 \pm 0.0100$	0.6547 ± 0.0037	$0.6483 \pm 0.0050$	$0.6422 \pm 0.0044$	0.6445 ± 0.0069
medical_charges_nominal	$0.8926 \pm 0.0002$	$0.8986 \pm 0.0002$	$0.8914 \pm 0.0003$	$0.8922 \pm 0.0001$	$0.8954 \pm 0.0021$	$0.8978 \pm 0.0011$	$0.8929 \pm 0.0002$	$0.8928 \pm 0.0001$	$0.8991 \pm 0.0007$	0.9033 ± 0.0019
megawatt_1	$0.8920 \pm 0.0099$	$0.8855 \pm 0.0068$	0.8802 ± 0.0098	$0.8907 \pm 0.0100$	$0.8920 \pm 0.0099$	0.8854 ± 0.0104	0.8973 ± 0.0069	0.8947 ± 0.0099	$0.9000 \pm 0.0022$	0.9039 ± 0.0099
messidor_features	$0.6539 \pm 0.0128$	$0.7197 \pm 0.0203$	$0.7219 \pm 0.0092$	$0.7416 \pm 0.0070$	$0.6733 \pm 0.0066$	$0.6794 \pm 0.0083$	$0.6814 \pm 0.0018$	0.6652 ± 0.0010	$0.6858 \pm 0.0196$	$0.6979 \pm 0.0071$
openml_586	$0.6619 \pm 0.0093$	$0.7162 \pm 0.0103$	$0.6666 \pm 0.0134$	$0.7108 \pm 0.0052$	$0.6880 \pm 0.0296$	$0.7735 \pm 0.0270$	$0.7200 \pm 0.0033$	$0.7200 \pm 0.0033$	$0.7185 \pm 0.0309$	0.7721 ± 0.0013
openml_589	$0.6557 \pm 0.0032$	$0.7022 \pm 0.0059$	$0.6750 \pm 0.0031$	$0.6870 \pm 0.0012$	$0.7208 \pm 0.0068$	$0.7711 \pm 0.0162$	$0.6961 \pm 0.0036$	$0.6961 \pm 0.0036$	$0.6954 \pm 0.0341$	$0.7351 \pm 0.0148$
openml_607	$0.6362 \pm 0.0083$	$0.7036 \pm 0.0005$	$0.6326 \pm 0.0092$	$0.6506 \pm 0.0126$	$0.6986 \pm 0.0535$	0.7655 ± 0.0081	$0.6916 \pm 0.0084$	$0.6916 \pm 0.0084$	$0.7225 \pm 0.0147$	0.7254 ± 0.0176
pima_indian	0.7391 ± 0.0007	$0.7574 \pm 0.0033$	0.7452 ± 0.0111	$0.7353 \pm 0.0133$	0.7404 ± 0.0041	0.7405 ± 0.0074	0.7587 ± 0.0138	$0.7505 \pm 0.0033$	0.7609 ± 0.0087	0.7543 ± 0.0098
poker_hand	0.6862 ± 0.0041	$0.6862 \pm 0.0041$	0.9973 ± 0.0001	$0.6862 \pm 0.0041$	1.0000 ± 0.0000	$1.0000 \pm 0.0000$	$0.7370 \pm 0.0365$	0.7856 ± 0.0579	$0.9965 \pm 0.0039$	1.0000 ± 0.0000
ps5_episode_3	$0.8511 \pm 0.0036$	$0.8490 \pm 0.0028$	0.8486 ± 0.0045	$0.8511 \pm 0.0036$	$0.8551 \pm 0.0013$	0.8560 ± 0.0018	$0.8601 \pm 0.0025$	$0.8565 \pm 0.0056$	0.8696 ± 0.0107	0.8554 ± 0.0019
ps5_episode_4	0.5767 ± 0.0003	$0.5791 \pm 0.0000$	0.5736 ± 0.0000	$0.5767 \pm 0.0003$	0.5771 ± 0.0001	0.5773 ± 0.0005	0.5772 ± 0.0001	$0.5772 \pm 0.0003$	$0.5775 \pm 0.0002$	0.5784 ± 0.0011
spectf	$0.7926 \pm 0.0057$	$0.7851 \pm 0.0265$	$0.7826 \pm 0.0247$	$0.7926 \pm 0.0057$	$0.7851 \pm 0.0021$	$0.8114 \pm 0.0056$	$0.8200 \pm 0.0100$	$0.8064 \pm 0.0076$	$0.8764 \pm 0.1007$	0.8278 ± 0.0210
statlog_german_credit	$0.7507 \pm 0.0038$	$0.7443 \pm 0.0080$	0.7457 ± 0.0087	$0.7523 \pm 0.0038$	$0.7490 \pm 0.0052$	$0.7533 \pm 0.0080$	$0.7547 \pm 0.0038$	$0.7597 \pm 0.0025$	$0.7587 \pm 0.0045$	$0.7683 \pm 0.0061$
wine_quality	$0.6575 \pm 0.0057$	$0.6600 \pm 0.0073$	$0.6545 \pm 0.0055$	$0.6564 \pm 0.0047$	$0.6545 \pm 0.0047$	$0.6575 \pm 0.0057$	$0.6623 \pm 0.0028$	$0.6619 \pm 0.0024$	$0.6572 \pm 0.0017$	$0.6608 \pm 0.0024$
Average Score	0.7558 ± 0.1017	0.7684 ± 0.0922	0.7718 ± 0.1101	0.7651 ± 0.0948	0.7791 ± 0.1068	0.7900 ± 0.0996	0.7745 ± 0.0958	0.7750 ± 0.0969	0.7908 ± 0.1105	0.7917 ± 0.1037

## A.3 EVALUATING PIFE UNDER REALISTIC TRAIN-TEST SPLITS

In our setup, we sought to closely mimic how feature engineering is typically performed in real-world applications by practitioners. Since a central objective of machine learning is to generalize effectively to unseen data, we designed our evaluation of PiFE to reflect this scenario. Each dataset was randomly split into training and test sets in a 70/30 ratio. The feature search process was restricted to the training split, where cross-validation was applied, while the test split was kept strictly unseen and used only for final evaluation. To mitigate order-related bias, all datasets were shuffled with fixed random seeds prior to splitting. This procedure may introduce slight variations in the reported scores compared to earlier tables. The detailed results of this evaluation are presented in Table 8.

Table 8: Performance comparison across competitions. Values are mean  $\pm$  std (tiny). Best values in each Train/Test column group are in bold.

Competition	Base	eline	CA	AFE	Ope	nFE	OC	Tree	PI	FE
Competition	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
adult	0.8493 ±0.0006	0.8495 ±0.0014	0.8504 ±0.0007	0.8487 ±0.0044	0.8530 ±0.0003	0.8420 ±0.0067	0.8526 ±0.0018	0.8518 ±0.0033	0.8504 ±0.0011	0.8517 ±0.0030
fertility	$0.8571 \pm 0.0286$	$0.9000 \pm 0.0333$	$0.8619 \pm 0.0360$	$0.8667 \pm 0.0333$	0.8571 ±0.0247	$0.8667 \pm 0.0000$	$0.8905 \pm 0.0218$	$0.8556 \pm 0.0192$	$0.9095 \pm 0.0360$	$0.8667 \pm 0.0000$
medical_charges_nominal	$0.8913 \pm 0.0006$	$0.8918 \pm 0.0004$	0.9001 ±0.0069	$0.8978 \pm 0.0049$	0.8975 ±0.0005	$0.8782 \pm 0.0058$	$0.8916 \pm 0.0004$	$0.8919 \pm 0.0002$	0.9000 ±0.0011	$0.8982 \pm 0.0027$
openml_586	$0.6221 \pm 0.0158$	$0.6603 \pm 0.0195$	0.7713 ±0.0260	$0.7461 \pm 0.0291$	0.6780 ±0.0327	$0.6992 \pm 0.0066$	$0.6927 \pm 0.0147$	$0.7133 \pm 0.0178$	0.7294 ±0.0170	$0.7271 \pm 0.0292$
pima_indian	$0.7344 \pm 0.0336$	$0.7287 \pm 0.0265$	0.7455 ±0.0244	$0.7359 \pm 0.0263$	0.7524 ±0.0195	$0.7460 \pm 0.0222$	0.7692 ±0.0153	$0.7273 \pm 0.0338$	$0.7698 \pm 0.0173$	$0.7677 \pm 0.0275$
ps5_episode_4	$0.5714 \pm 0.0006$	$0.5753  \pm 0.0008$	0.5738 ±0.0027	0.5772 ±0.0030	0.5744 ±0.0002	$0.5557 \pm 0.0005$	0.5721 ±0.0001	$0.5757 \pm 0.0006$	0.5733 ±0.0004	$\overline{0.5703}_{\pm 0.0039}$
Average	0.7543 ±0.1268	0.7676 ±0.1268	0.7838 ±0.1123	0.7787 ±0.1128	0.7688 ±0.1182	0.7646 ±0.1173	0.7781 ±0.1199	0.7693 ±0.1134	0.7887 ±0.1205	0.7803 ±0.1148

The results indicate that PIFE consistently delivers strong performance, achieving both high cross-validation scores on the training data and correspondingly high accuracy on the test data across a large majority of datasets. This alignment suggests that effective feature engineering plans discovered by PIFE are not only tuned for the training split but also generalize reliably to unseen data, reinforcing the robustness of the approach.

#### A.4 REVERSE POLISH NOTATION FOR FEATURE REPRESENTATION

We adopt Reverse Polish Notation (RPN) from Zou et al. (2025) as a representation scheme for the features generated in PIFE. In RPN, operators follow their operands, eliminating the need for parentheses and reducing ambiguity in expression evaluation. This structure allows for a compact and unambiguous encoding of complex feature transformations, which is particularly useful when features are generated programmatically or by language models.

Using RPN provides several advantages. First, it enables straightforward reconstruction of the original feature expression, as the sequence of operands and operators directly encodes the computational order. Second, RPN facilitates efficient storage and manipulation of features, since it can be easily parsed into computational graphs or evaluated using stack-based execution.

#### A.5 HYPERPARAMETERS

#### A.5.1 PREDICTIVE MODELS

As shown in Table 9, we consider both regression and classification models with standard hyperparameters for baseline evaluation and testing AutoFE methods.

Table 9: Regression and Classification Models with Parameters

Task	Model	Parameters
Regression	Linear Regression	_
Regression	Random Forest Regressor	n_estimators = 10, random_state = 0
Regression	XGBoost Regressor	n_estimators = 10, random_state = 0
Classification	Logistic Regression	solver = saga, class_weight = balanced, tol =
		$0.0005$ , C = $0.5$ , max_iter = $10000$ , penalty = $12$
Classification	Random Forest Classifier	n_estimators = 10, random_state = 0
Classification	XGBoost Classifier	n_estimators = 10, random_state = 0

#### A.5.2 PARAMETERS FOR DEEP LEARNING METHODS

We list down hyperparameters used for MLP and HyperFast in Table 10 and Table 11.

Table 10: Parameters for MLP

Parameter	Distribution / Range
Number of layers	UniformInt [1, 8]
Layer size	UniformInt [16, 1024]
Dropout	Uniform [0, 0.5]
Learning rate	LogUniform $[1 \times 10^{-5}, 1 \times 10^{-2}]$
Category embedding size	UniformInt [64, 512]
Learning rate scheduler	[True, False]
Batch size	[256, 512, 1024]

Table 11: Parameters for HyperFast

Parameter	Distribution / Range
Number of ensembles (N)	[1, 4, 8, 16, 32]
Batch size	[1024, 2048]
NN bias	[True, False]
Stratified sampling	[True, False]
Optimization strategy	[None, optimize, ensemble_optimize]
Optimize steps	[1, 4, 8, 16, 32, 64, 128]
Random seed	UniformInt [0, 9]

## A.5.3 FEATURE SELECTION METHODS

Details of parameters used for feature selection experiments (results in Table 3) are listed in Table 12 and Table 13.

Table 12: Parameters for Bayesian CMI-based Feature Selection

Parameter	Values
alpha	0.5
trials	50
distance	gower
min_num_feat_selected	0
scaling_criteria	min_max
sample_df	True
k	10
denomination	2
num_select_features	3

## A.5.4 AUTOFE METHODS

Table 14 summarizes the key parameters and their values for the AutoFE methods evaluated in this study, including PIFE, OCTREE, CAAFE, OPENFE, AUTOFEAT, and DFS. These values were chosen based on prior literature and preliminary experiments to ensure fair and comparable evaluation across methods.

## A.6 BASELINE SELECTION

While learning-based approaches, which aim to learn transformation policies directly from data, represent an important category, we do not include them in our current evaluations due to the high implementation complexity and substantial computational cost involved in training and adapting these models. This selection allows us to compare how different strategies perform in practice and to analyze their respective strengths, limitations, and implications for the future of feature engineering.

Table 13: Parameters for Genetic Algorithm-based Feature Selection

Parameter	Values
elitism	5
generations	30
population_size	30
crossover_prob	0.8
mutation_prob	0.05

Table 14: Parameters and Values for AutoFE Variants

Method	Parameter	Value
PIFE	total_steps	10
	eda_steps	3
	max_plots	3
	max_insights_per_plot	1
	n_features	5
	timeout	86400
OCTree	total_steps	5
	rule_steps	10
	n_features	1
	timeout	86400
CAAFE	total_steps	10
	n_features	1
	n_repeats	1
	timeout	86400
OpenFE	min_candidate_features	2000
	feature_boosting	False
	n_repeats	1
	timeout	86400
Autofeat	feateng_steps	2
	timeout	86400
DFS	max_depth	2
	transformations	transform_primitives
	timeout	86400

The original CAAFE framework was primarily designed for classification tasks and evaluated only with accuracy as the performance metric. In our adaptation, we extend CAAFE to also support regression tasks, introduce additional evaluation metrics beyond accuracy for a fairer comparison, and enrich the set of operators available for feature construction.

OCTree, on the other hand, required more substantial modifications. The original implementation was tightly coupled with specific LLM APIs and lacked iterative refinement. We restructured its pipeline to generalize API usage, extended the feature generation loop to be iterative, and incorporated support for regression tasks along with additional evaluation metrics. These modifications make OCTree more robust and applicable across a broader range of tabular learning scenarios.

#### A.7 DISCUSSION ON FEATURE SELECTION METHODOLOGIES

Feature selection in our framework can be applied at two stages: (1) immediately after the feature engineering step, or (2) after the full PiFE run. In our experiments, Feature Importance and Bayesian Conditional Mutual Information (CMI) based selection were applied after the feature engineering stage but prior to model validation, whereas a Genetic Algorithm based selection was performed after the complete pipeline execution.

Feature selection plays a crucial role in enhancing both model interpretability and generalization. While an individual engineered feature may appear weak in isolation, its combination with other features can capture complex interactions and yield a much stronger predictive signal. Without an appropriate selection mechanism, such subtle but useful interactions may be overlooked or drowned out by a large number of irrelevant or redundant features. By systematically ranking and filtering features, our selection strategies help retain those that contribute jointly to predictive power, thereby improving efficiency, reducing overfitting, and uncovering more meaningful feature representations.

#### A.7.1 FEATURE IMPORTANCE BASED SELECTION

After features are generated in an iteration , we perform validation on the dataset with new features and compute feature importance scores. If the validation score of the current run is greater than that of the previous run, we retain all the features. Otherwise, we apply a filtering criterion: only those features with importance greater than 1 / (number\_of\_features) are selected. This threshold is motivated by the expectation that a retained feature should contribute at least more than the average share of importance across all features.

#### A.7.2 CMI-BASED BAYESIAN OPTIMIZATION FOR FEATURE GROUP SELECTION

**Bayesian Optimization (BO).** BO is a sequential strategy for optimizing expensive black-box functions. A surrogate model (e.g., Gaussian Process) provides mean  $\mu(x)$  and uncertainty  $\sigma(x)$ , guiding the selection of new points via an acquisition function a(x) (e.g., EI, UCB):

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} a(x \mid \mu(x), \sigma(x)).$$

In feature engineering, we need to add a feature subset to the existing feature set. BO treats this feature subset as X and the CMI as f(X)(objective function), enabling efficient exploration of feature combinations.

$$f(X) = I(X; Y \mid Z),$$

where X is the feature subset, Y is the target, and Z is the feature set.

Conditional Mutual Information (CMI). CMI in feature engineering quantifies the unique contribution of a feature subset X to predicting Y given the feature set Z:

$$I(X;Y\mid Z) = \iiint p(x,y,z) \log \frac{p(x,y\mid z)}{p(x\mid z) p(y\mid z)} dx dy dz.$$

We use a slightly modified version of the mixed-type k-NN estimator from Mesner & Shalizi (2020), which is robust to discrete and continuous variables. We set  $K = \max(3, \min(20, \sqrt{n}))$ , which helps mitigatethe high-dimensionality issue in CMI calculation.

#### A.7.3 GENETIC ALGORITHM

Genetic Algorithms (GAs) are population-based metaheuristic optimization methods. A GA maintains a population of candidate solutions (chromosomes), where each chromosome encodes a subset of features (typically as a binary string, with 1 indicating selection of a feature and 0 otherwise). The algorithm evolves this population through the iterative application of genetic operators:

- Selection: Chromosomes are chosen based on their fitness, which in our case is the predictive performance (e.g., validation accuracy or F<sub>1</sub> score)
- Crossover: Pairs of chromosomes exchange parts of their feature subsets, enabling exploration of new feature combinations.
- Mutation: Random flips of feature bits introduce diversity and help escape local optima.

The fitness of a chromosome c can be expressed as

$$Fitness(c) = Score(f(X_c), Y),$$

where  $X_c$  denotes the features selected by chromosome c, Y is the target, and  $f(\cdot)$  is the downstream predictive model.

Over successive generations, the GA converges toward feature subsets that maximize predictive performance. While computationally more expensive than other methods like CMI-BO and MFI, GAs often identify subsets of features with strong predictive power, making them effective when interactions between features play an important role.

## A.8 PROMPTS

864

865

866

#### Listing 1: EDA Analysis Code Generation

```
867
       You are an EDA agent operating in a Kaggle Grandmaster-style automated
868
          feature engineering pipeline. This is iteration {current_iteration}
869
          of a multi-stage EDA loop. Your role is to produce competition-grade
870
          exploratory data analysis code that progressively builds upon the
871
          analyses performed in previous iterations.
872
873
       Remember: this is a strategic, hypothesis-driven EDA process - think
          like a top Kaggle competitor uncovering hidden signal iteratively
874
          across multiple passes.
875
876
      You are provided with:
877
       - Dataset description
       - Summary of preprocessing steps taken
878
       - Previous EDA code history and respective observations if any.
879
880
       Dataset Description:
881
      {dataset_description}
882
       Pre-processing Steps:
883
       {preprocessing_summary}
884
885
      Focus Areas:
886
       {focus areas}
887
      Memory (previous code and observations):
888
       {memory}
889
890
891
892
       ### Analysis Constraints
      STRICT LIMITS:
893
       - Maximum {max_plots} plots total
894
       - Do NOT explicitly suggest feature transformations, binning, encoding,
895
          or normalization.
896
       - Focus on uncovering patterns, trends, correlations, and anomalies in
897
          the data.
       - Avoid bias towards only high-importance features include a mix of
898
          numerical, categorical, and temporal features.
899
       - For large datasets (>5000 rows), sample strategically before complex
900
          plots.
901
902
903
       ### Response Format
904
      Your response should strictly follow the following Code Structure:
905
906
      Before the code block include a short **Implementation Rationale** of 24
907
          sentences that explains:
       - Why you chose the specific analyses / plots (what hypothesis you are
908
          testing),
909
       - What you expect the output to reveal (the type of insight sought),
910
       - One-line failure mode / limitation: Think Harder (e.g., 'may fail on
911
          heavy-tailed column; will downsample if >5000 rows').
912
       - Divide the code into code cells using '# %%' to demarcate different
913
          sections in the code.
914
       - Code should be wrapped within '''python ... ''' quotes. Do not write
915
          code at any other place than this.
916
       - Do not write the context of the code block in the same line as '# %%'.
917
          Write it in a new line with enumeration, where enumeration should be
          in comments.
```

```
918
       - Each section should focus on a different type of analysis aimed at
919
          revealing meaningful patterns.
920
       - Ensure diversity by analyzing features across different types and
921
          varying levels of correlation with the target.
       - All plots must have proper titles, axis labels, and legends where
922
          applicable.
923
       - Produce plots that are clear and informative, suitable for
924
          presentation.
925
       - If plots become too crowded or contain too many elements to be
926
          readable, split them into multiple smaller plots.
       - Each plot should be individually assigned to a unique and
927
          human-readable variable name.
928
       - Do NOT use 'plt.show()' or save figures to files just generate them.
929
       - For expensive plots (swarm, violin, scatter, kde, displot, etc.), if
930
          dataset size >5000 rows, downsample to 5000 using stratified
931
          sampling (if categorical column available), else use appropriate
          sampling.
932
933
934
935
       ### Mandatory Imports
936
       Each EDA code block must begin with:
       '''python
937
       import pandas as pd
938
       import numpy as np
939
       import matplotlib.pyplot as plt
940
       import seaborn as sns
941
942
943
       ### General Instructions
944
945
       - Assume the dataset with variable name df is already present.
946
       - Always include all required imports in your code. Do NOT assume
          imports persist across steps.
947
       - Preserve existing variable assignments and do not overwrite previously
948
          assigned variables.
949
       - Maintain continuity with past EDA steps by avoiding duplicate analyses
950
          and expanding on previously explored patterns.
951
       - Incorporate findings from earlier iterations to guide where to dig
          deeper.
952
       - Ensure coverage across diverse feature types (numerical, categorical,
953
          temporal).
954
       - Generate new, non-redundant insights that add incremental value.
955
       - All generated code should be clean, modular, and ready for execution
956
          without edits.
       - Use # %% to clearly separate analysis sections.
957
       - Include explanatory comments in the format: # INSIGHT: <purpose of
958
          this analysis>, but do NOT explicitly suggest feature
959
          transformations.
960
```

#### Listing 2: EDA Analysis Insight Generation

961

962 963

964

965

966

967

968 969

970

```
You are a feature engineering specialist analyzing outputs from an EDA process.
You are provided with:

- The code used to generate the EDA

- The plots generated from that code

- The textual/statistical outputs produced

Your job is to extract meaningful, high-value insights from these materials and then propose specific, well reasoned feature transformations inspired by these insights. Insights must be grounded in visual and statistical evidence, not speculation.
```

```
972
973
      Think like a Kaggle Grandmaster preparing features for a top-tier
974
          competition.
975
       # Task Description
976
       {dataset_description}
977
978
       # EDA Code Context
979
       {eda_code}
980
       # Available Operators
981
       {operators_description}
982
983
       # Instructions
984
985
       ## Analysis Guidelines When forming INSIGHTS:
       - Highlight relationships, anomalies, patterns, or distributions that
986
          stand out in the data.
987
       - Capture interactions, trends, and category-level differences.
988
       - Refer to the specific visualization or statistical summary they come
989
990
       - Generate maximum {max_insights_per_plot} insights per plot.
991
       ## When forming FEATURE_TRANSFORMATIONS:
992
       - Map each transformation to a corresponding insight.
993
       - Use ONLY the Available Operators listed above for suggesting feature
994
          transformations.
       - You can suggest dropping features if they are redundant, highly
995
          correlated, or unhelpful for the target variable.
996
       - Provide reasoning linked to potential model performance improvements.
997
       - Keep recommendations actionable, clear, and technically precise.
998
       - Strictly NEVER include the target column in any feature engineering,
999
          transformations, encodings, interaction terms, binning, scaling, or
          statistical computations.
1000
1001
1002
       # Response format
1003
1004
      Your output MUST be structured in exactly two sections using the
1005
          following XML-style tags:
1006
       <INSIGHTS>
1007
      List clear, evidence-backed observations from the EDA results and plots.
1008
      Avoid feature suggestions here keep this purely as descriptive,
1009
          analytical findings.
      Each insight should explicitly reference the plot or analysis it came
1010
          from.
1011
       </INSIGHTS>
1012
1013
      <FEATURE_TRANSFORMATIONS>
1014
      For each transformation:
1015
       - Specify the exact feature(s) involved
       - Describe the suggested transformation or engineering step using
1016
          Available Operators
1017
       - Provide a short reasoning for why it is beneficial based on the
1018
          insights above
1019
       - Include 35 high-priority transformations that would add the most value
       - You can suggest dropping features if they are redundant, highly
1020
          correlated, or unhelpful for the target variable.
1021
       </FEATURE_TRANSFORMATIONS>
1022
1023
      AVOID: generic or obvious patterns, restating axis labels, or vague
1024
          statements.
      FOCUS: insights that directly inform strong, competition-grade feature
1025
          engineering.
```

1026 1027 1028 Listing 3: Feature Engineering Code Generation 1029 1030 You are a feature-engineering agent. Your goal is to generate new 1031 machine-learning-ready features and return executable Python code 1032 that creates them. You will be provided with the Dataset Description, 1033 the Pre-processing Steps, the Feature Transformation Guidelines, 1034 list of allowed operators for feature engineering, the Feature importance scores and the Rejected Features. 1035 1036 # Task description: 1037 {dataset\_description} 1038 1039 # Pre-processing Steps: {preprocessing\_step\_summary(preprocess, target\_encoder)} 1040 1041 # Feature Transformation Guidelines: 1042 {quidelines} 1043 1044 # Available Operators: {operators\_description} 1045 1046 # Feature importance scores: 1047 {[f"{k}: {v:.5f}" for k, v in feature\_importance\_scores.items()] if 1048 feature\_importance\_scores is not None else []} 1049 # Rejected Features: 1050 {rejected\_features} 1051 1052 # Instructions: 1053 1054 ## Feature Engineering Instructions: 1) Generate exactly {n\_features} new machine-learning-ready features. 1055 2) You should think about the reasons for the rejected features and try 1056 to incorporate that learning while creating new featurs 1057 3) Feature Transformation Guidelines contains ideas based on exploratory 1058 data analysis conducted earlier. You should create new features that 1059 are based on the ideas in the Feature Transformation Guidelines. 4) You should use the list of Available Operators, for feature 1060 engineering to create new features. 1061 1062 ## Code Instructions: 1063 - Do not wrap the entire code inside a function or class. 1064 - Assume the environment is similar to a Jupyter Notebook, so you may use # %% to separate code blocks. You may define small utility/helper 1065 functions if needed, but make sure they are invoked within the same 1066 code block. 1067 - The final output should be an executable code block, not a function or 1068 class definition. 1069 - Ensure that code is enclosed within python code literal as follows. Do not write code anywhere else. 1070 '''python 1071 [code goes here] 1072

# Target Leakage Prevention Rules: STRICT RULES TO AVOID TARGET LEAKAGE: - NEVER create features that directly transform or encode the target column itself (e.g., log(target), residuals vs. target, deviation-from-target, ranks of target within bins, z-scores, percentiles, etc.). 20

1073 1074

1075

1076 1077

1078

```
1080
      - The target column may only be used to compute group-level aggregate
1081
          statistics (mean, median, std, count) at the group level.
1082
          - Example: mean(target) by year, median(target) by region.
           - These must be computed on df_train only, stored as a mapping, and
1083
          applied to df_test with a fallback.
1084
       - Forbidden feature patterns include:
1085
          - Any function that directly transforms target values row-by-row
1086
          (log, rank, residual).
1087
           - Any feature whose definition requires subtracting or dividing a
1088
          rows own target from an aggregate.
           - Any within-bin or within-group ranking of target values.
1089
          - Always check that the engineered feature can be computed in
1090
          exactly the same way for both df_train and df_test without needing
1091
          df_test[TARGET_COL].
1092
          - If a proposed feature would violate these rules, DO NOT generate
          it instead, list it in
1093
          SAFETY_REPORT['features_dropped_due_to_no_test_support'].
1094
1095
       # Feature Inclusion & Target-Use Rules (MANDATORY, concise):
1096
      CRITICAL: The runtime convention is that the target column WILL be
1097
          present in df_train as 'target', and WILL NOT be present in df_test.
1098
          Follow these rules without exception:
1099
       - INPUT ASSUMPTION
1100
          - df_train contains the target column named 'target'.
1101
           - df_test, if provided, MUST NOT contain the target column. Agents
1102
          must treat df_test as unlabeled.
1103
      - FEATURE INCLUSION: Every suggested feature must be constructible on
          BOTH df_train and df_test. If the feature cannot be created for
1104
          df_test without reading the target (TARGET_COL) or other unavailable
1105
          test-only data, DO NOT create that feature for df_train drop it
1106
          entirely. Never produce train-only features.
1107
      - TARGET USAGE (TRAIN-ONLY STATISTICS): You may compute aggregate
1108
          statistics using df_train[TARGET_COL] only to build train-derived
1109
          mapping objects (e.g., group means/counts/medians). All such
1110
          computations MUST:
1111
          - be computed only on df_train,
1112
      - APPLYING MAPPINGS TO TEST: For every train-derived mapping, provide
          explicit code that applies the mapping to df_test inside if df_test
1113
          is not None: using map/merge and .fillna(<fallback>).
1114
           - Example pattern (must be used):
1115
           '''python
1116
          mapping = df_train.groupby('X')[TARGET_COL].median().to_dict()
1117
          df_train['f'] = df_train['X'].map(mapping).fillna(<fallback>)
1118
          if df_test is not None:
               df_test['f'] = df_test['X'].map(mapping).fillna(<fallback>)
1119
          _train_mappings['mapping_name'] = mapping
1120
1121
       - INTERMEDIATE / TRAIN-ONLY COLUMNS: If you create intermediate columns
1122
          on df_train solely to compute mapping/statistics or to use them to
          create suggested features, remove them from df_train before
1123
          finishing the code (so columns remain symmetric). Do NOT leave
1124
          intermediate columns that cannot be created on df_test.
1125
      - COLUMN SYMMETRY CHECK: At the end of your code, ensure df_train and
1126
          df_test have the same columns (except for TARGET_COL in df_train).
1127
          Add this assertion:
1128
           '''python
           # Ensure column symmetry between train and test sets
1129
           if df_test is not None:
1130
              train_cols = set(df_train.columns) - {TARGET_COL}
1131
              test_cols = set(df_test.columns)
1132
               assert train_cols == test_cols, f"Column mismatch: train has
          {train_cols - test_cols} extra, test has {test_cols - train_cols}
1133
          extra"
```

```
1134
1135
      - NON-IMPLEMENTABLE FEATURES: If any transform (e.g., direct arithmetic
1136
          with TARGET_COL in df_test, or features requiring target values at
1137
          test-time) cannot be implemented safely on df_test, explicitly
          exclude that feature and list it under
1138
          SAFETY_REPORT['features_dropped_due_to_no_test_support'].
1139
       - FORBIDDEN: Under no circumstance should the generated code access
1140
          TARGET_COL within df_test or assume its existence there. Do not
1141
          create features that would require test-time predictions or labels.
1142
      - RANDOMNESS: Use deterministic randomness via RANDOM_STATE for any
          sampling/splitting operations on df_train; do not sample df_test.
1143
1144
      # Templates (must be used for any transform that depends on target/train
1145
          statistics):
1146
      Provide transformations using these exact patterns when the operation
          depends on train statistics.
1147
1148
      A) GroupByThenMean (safe pattern)
1149
       '''python
1150
      # compute mapping on train ONLY
1151
      edu_mean_by_occ =
          df_train.groupby('occupation')['education-num'].mean().to_dict()
1152
       # apply to train
1153
      df_train['QualificationSurplus'] = df_train['education-num'] -
1154
          df_train['occupation'].map(edu_mean_by_occ)
1155
       # apply to test (no target used). fallback to 0 for unseen occupations
1156
      if df_test is not None:
1157
          df_test['QualificationSurplus'] = (
               df_test['education-num']
1158
          df_test['occupation'].map(edu_mean_by_occ)
1159
          ).fillna(0)
1160
1161
      B) Target-like encoding (train-derived, safe pattern)
1162
       '''python
1163
       # compute target-encoding stats on train ONLY
1164
      enc_by_cat = df_train.groupby('cat_col')[TARGET_COL].agg(
1165
                       ['mean','count']
1166
                   ).to_dict(orient='index')
      # convert to mapping (use mean, with global fallback)
1167
      global_mean = df_train[TARGET_COL].mean()
1168
      cat_mean_map = {k: v['mean'] for k, v in enc_by_cat.items()}
1169
       # apply
1170
      df_train['cat_col_te'] =
1171
          df_train['cat_col'].map(cat_mean_map).fillna(global_mean)
      if df_test is not None:
1172
          df_test['cat_col_te'] =
1173
          df_test['cat_col'].map(cat_mean_map).fillna(global_mean)
1174
1175
1176
      C) Stratified sampling for train-only operations (must not touch df_test)
1177
       '''python
1178
      from sklearn.model_selection import StratifiedKFold
1179
      skf = StratifiedKFold(n_splits=5, shuffle=True,
1180
          random_state=RANDOM_STATE)
1181
      for train_idx, holdout_idx in skf.split(df_train, df_train[TARGET_COL]):
           # operate only on df_train.iloc[train_idx], use holdout for internal
1182
          validation
1183
          pass
1184
      . . .
1185
1186
      D) ALWAYS include mapping objects in code and show how they'll be
1187
          persisted/serialized if needed.
```

```
1188
      # Templates for Intermediate Features:
1189
1190
       '''python
1191
      Interaction: lymphatics early_uptake (concatenated string, factorized)
      df_train['lymphatics_earlyuptake'] = df_train['lymphatics'].astype(str)
1192
          + '_' + df_train['early_uptake'].astype(str)
1193
      # Factorize (shared mapping for train, then reapply to test)
1194
      all_cats = pd.concat([df_train['lymphatics_earlyuptake'],
1195
                           (df_test['lymphatics'].astype(str) + '.
          df_test['early_uptake'].astype(str)) if ('df_test' in locals() and
1196
          df_test is not None) else pd.Series([], dtype=str)]
1197
      lympt_early_map, lympt_early_uniques = pd.factorize(all_cats, sort=True)
1198
      train_codes = lympt_early_map[:len(df_train)]
1199
      df_train['lymphatics_earlyuptake_code'] = train_codes
1200
      # Remove lymphatics_earlyuptake as its not the suggested feature and it
          can not be created in df_test
1201
      df_train.drop(columns=['lymphatics_earlyuptake'], inplace=True)
1202
      if 'df_test' in locals() and df_test is not None:
1203
           test_codes = lympt_early_map[len(df_train):]
1204
           df_test['lymphatics_earlyuptake_code'] = test_codes
1205
       _train_mappings['lymphatics_earlyuptake_factorization'] =
1206
          dict(zip(lympt_early_uniques, range(len(lympt_early_uniques))))
1207
1208
       - IMPORTANT: Intermediate Feature Cleanup
1209
      At the end of your code, ensure you remove ALL intermediate features
1210
          that were created solely for computation purposes:
1211
1212
       # Clean up intermediate features
1213
      intermediate_features = ['temp_feature1', 'temp_feature2',
1214
          'mapping_temp']
1215
      for feature in intermediate_features:
           if feature in df_train.columns:
1216
               df_train.drop(columns=[feature], inplace=True)
1217
               if df_test is not None and feature in df_test.columns:
1218
                   df_test.drop(columns=[feature], inplace=True)
1219
       . . .
1220
1221
       # Response Format for Python Code:
1222
       - Python code for n feature transformations
1223
1224
       '''python
1225
       [feature engineering code]
1226
1227
       - RPN Format:
1228
1229
      Reverse Polish Notation and Description of n feature transformations
1230
      Provide in this format:
1231
      FeatureName: <new_feature_name>
1232
      RPN : feature1 feature2 +
1233
      Description : Sum of feature1 and feature2
1234
1235
      FeatureName: <new_feature_name>
1236
      RPN: feature3 feature4 GroupByThenMean -
      Description: Difference between feature3 and mean of feature3 grouped
1237
          by feature4
1238
1239
       # Instructions for RPN Notation:
1240
1241
      - Use the format Dropped_<FeatureName> for features that are dropped.
      - Do not use square brackets in the FeatureName.
```

```
1242
      - Separate tokens in the RPN string with spaces.
1243
      - Examples of correct RPN:
1244
          - feature1 feature2 +
           - feature1 drop (for dropping a feature)
1245
           - feature1 feature2 GroupByThenMean -
1246
1247
       - Avoid invalid RPN such as featurel feature2 GroupByThenMean -
1248
          subtraction - requires two operands.
1249
       - Example with GroupByThenMean:
1250
1251
           '''python
1252
          edu_mean_by_occ =
1253
          df_train.groupby('occupation')['education_num'].mean()
1254
          df_train['QualificationSurplus'] = df_train['education-num'] -
          df_train['occupation'].map(edu_mean_by_occ)
1255
          python
1256
          Copy code
1257
          edu_mean_by_occ :
1258
          df_train.groupby('occupation')['education_num'].mean()
          df_train['QualificationSurplus'] = df_train['education-num'] -
1259
          df_train['occupation'].map(edu_mean_by_occ)
1260
           if df_test is not None:
1261
               df_test['QualificationSurplus'] = (
1262
                   df_test['education-num'] -
1263
          df_test['occupation'].map(edu_mean_by_occ)
1264
              ).fillna(0)
1265
1266
          RPN: education-num education_num occupation GroupByThenMean -
1267
          This RPN correctly represents the operation.
1268
           Incorrect RPN: education_num occupation GroupByThenMean - (only 1
1269
          operand before subtraction)
1270
       - Drop Operation Examples:
1271
           - To drop a feature: RPN: feature_name drop
1272
           - To drop multiple features: RPN: feature1 drop feature2 drop
1273
           - Always document dropped features in your response with the
1274
          Dropped_<FeatureName> format
1275
```

## A.9 Example PIFE Features from experiment runs

1276

1277

1278 1279

1280

1281

1282

1283

1284

1285

1286

1287

1289

1290

1291

1292

```
Higher Order Feature
Competition: spectf
Name of the feature: Sum_x_Hotspot
RPN: F1S F2S + F3S + F4S + F5S + F6S + F7S + F8S + F9S + F10S + F11S + F12S +
F13S + F14S + F15S + F16S + F17S + F18S + F19S + F20S + F21S + F22S + F20S F21S
max F22S F13S max max *
Order: 21
EDA Reasoning: Interaction of global stress and regional hotspot (maps to the
"global-regional synergy" insight)
Features: F1S,...,F22S, max_stress_13_20_21_22
Transformation:
1) StressSum_all = F1S + F2S + ... + F22S (chain the "+" operator across all stress ROIs)
2) Sum_x_Hotspot = StressSum_all max_stress_13_20_21_22
Reasoning: The Q1 quadrant (high-high) showed a 0.43 abnormal rate vs 0.20 elsewhere.
The multiplicative term encodes this synergy explicitly and is often more predictive than
either marginal.
```

Figure 5: Higher Order Feature Generation

#### **Higher Predictive Power**

Competition: poker\_hand

Name of the feature: rank\_pair\_sum

**RPN**: C1 C1 / C1 C2 - abs + reciprocal round C1 C1 / C1 C3 - abs + reciprocal round + C1 C1 / C1 C4 - abs + reciprocal round + C1 C1 / C1 C5 - abs + reciprocal round + C2 C1 - abs C1 C1 / + reciprocal round + C2 C3 - abs C1 C1 / + reciprocal round + C2 C4 - abs C1 C1 / + reciprocal round + C2 C5 - abs C1 C1 / + reciprocal round + C3 C4 - abs C1 C1 / + reciprocal round + C3 C5 - abs C1 C1 / + reciprocal round + C4 C5 - abs C1 C1 / + reciprocal

**Feature Importance**: 0.35005184128253863 Increase in score after adding feature: 0.306

**EDA Reasoning**: Feature(s): C1–C5

**Transformation**: Build pairwise "same-rank" indicators and aggregate. 1) ONES = C1 / C1 2) For each unordered pair (i, j) among 1..5: diff\_ij = abs(Ci Cj) diff1\_ij = diff\_ij + ONES eq\_ij = round(reciprocal(diff1\_ij)) equals 1 if Ci=Cj, else 0 3) rank\_pair\_sum = sum(eq\_ij over the 10 pairs) using + 4) For each i: eq\_i = sum(eq\_ij over j i) using + max\_same\_rank = max(eq\_1, eq\_2, eq\_3, eq\_4, eq\_5) using max

**Reasoning**: From the correlation heatmap and uniform marginals, single ranks are uninformative; equality patterns drive CLASS. rank\_pair\_sum differentiates high-card/straight/flush (0), one pair (1), two pair (2), three-kind (3), full house (4), four-kind (6). max\_same\_rank (values 0–3) is a strong proxy for the largest multiplicity (pair/three/four). These directly target rare classes (3,6,7) and improve separability under heavy imbalance.

## Figure 6: Higher Predictive Power

#### Weak features combined to get strong feature

**Competition**: messidor\_features **Name of the feature**: MA\_RANGE

**RPN**: ma1 ma2 max ma3 max ma4 max ma5 max ma6 max ma1 ma2 min ma3 min ma4 min ma5 min ma6 min -

## **Feature Importances:**

MA\_RANGE: 0.09319
ma1: 0.04985
ma2: 0.04468
ma3: 0.04219
ma4: 0.03743
ma5: 0.04279
ma6: 0.03985

**EDA Reasoning**: MA\_RANGE: max(ma1, max(ma2, max(ma3, max(ma4, max(ma5, ma6))))) min(ma1, min(ma2, min(ma3, min(ma4, min(ma5, ma6)))))

**Reasoning**: Despite high within-cluster correlations, thresholds are not identical. The boxplots show distributional spread growing with DR. The range across thresholds measures stability/sensitivity of detections to  $\alpha$ ; noisy/non-DR images may show different spreads than true DR.

Figure 7: Effective Feature Combination

#### A.10 LLM USAGE

Apart from our proposed framework, PiFE, we leveraged LLMs to assist in refining the writing of this research paper. The models were used solely for language polishing, grammar corrections, and clarity improvements, without influencing the scientific content, experimental design, results, or conclusions.

#### A.11 Broader Impact Statement

**Utility and Real-World Relevance.** The AutoFE method can autonomously generate, evaluate, and select meaningful features from raw data, potentially enabling more robust and interpretable predictive models in real-world applications. By leveraging EDA-driven insights, it can provide data scientists with explainable and interpretable features, enhancing model transparency and decision-making. With appropriate statistical safeguards, validation checks, and privacy-preserving measures, it can help mitigate the risk of spurious or misleading features and support the responsible deployment of LLM-assisted feature engineering systems.

Risks and Biases. LLM-assisted feature generation can be influenced by biases present in the training data, including memorization of datasets or solutions from prior competitions, which may lead to overfitting or inflated performance on familiar tasks. To mitigate this, incorporating recent Kaggle competitions and unseen datasets during evaluation can help assess generalization and reduce reliance on memorized patterns. Selecting diverse datasets from multiple sources is critical to capture varied real-world scenarios, minimize systemic bias, and ensure broadly applicable and fair features. Additionally, running experiments across multiple random seeds can provide a more robust assessment of LLM-based frameworks, helping to quantify variability and improve reliability in feature generation.