

Knowledge-Informed Automatic Feature Extraction via Collaborative Large Language Model Agents

Henrik Brådlund^{1 2 3} Morten Goodwin² Vladimir I. Zadorozhny^{1 2} Per-Arne Andersen²

Abstract

The performance of machine learning models on tabular data is critically dependent on high-quality feature engineering. While Large Language Models (LLMs) have shown promise in automating feature extraction (AutoFE), existing methods are often limited by monolithic LLM architectures, simplistic quantitative feedback, and a failure to systematically integrate external domain knowledge. This paper introduces Rogue One, a novel, LLM-based multi-agent framework for knowledge-informed automatic feature extraction. Rogue One operationalizes a decentralized system of three specialized agents—Scientist, Extractor, and Tester—that collaborate iteratively to discover, generate, and validate predictive features. Crucially, the framework moves beyond primitive accuracy scores by introducing a rich, qualitative feedback mechanism and a "flooding-pruning" strategy, allowing it to dynamically balance feature exploration and exploitation. By actively incorporating external knowledge via an integrated retrieval-augmented (RAG) system, Rogue One generates features that are not only statistically powerful but also semantically meaningful and interpretable. We demonstrate that Rogue One significantly outperforms state-of-the-art methods on a comprehensive suite of 19 classification and 9 regression datasets. Furthermore, we show qualitatively that the system surfaces novel, testable hypotheses, such as identifying a new potential biomarker in the myocardial dataset, underscoring its utility as a tool for scientific discovery.

¹School of Computing and Information, University of Pittsburgh, Pittsburgh, USA ²Centre for Artificial Intelligence Research, University of Agder, Grimstad, Norway ³Norkart AS, Oslo, Norway. Correspondence to: Henrik Brådlund <henrik.bradland@uia.no>.

1. Introduction

The rising need for interpretable machine learning, combined with their documented challenges on smaller tabular datasets (Grinsztajn et al.), has led to the continued dominance of classical tree-based models like Random Forest and XGBoost for tabular prediction. Nevertheless, the performance of the tree-based models is heavily tied to the predictive power of the underlying features. This introduces a shift from a model-centric development to a data-centric development, where synthesizing strong features is the driver for stronger downstream performance. Traditional approaches use the intuition of domain experts to generate strong features. Although effective, this method is highly subjective and does not well explore the space of possible features. Research has proposed to automate the process, thus giving rise to the field of Automatic Feature Extraction (AutoFE). Traditional AutoFE methods (e.g., polynomial feature generation) struggled to create semantically meaningful features, as they lacked the ability to incorporate domain knowledge or understand the high-level relationships between feature concepts. Nevertheless, with the rise of Large Language Models (LLMs), with their strong contextual understanding and parametric knowledge, they represent a powerful new paradigm for discovering new features as demonstrated by Abhyankar et. al. (Abhyankar et al., 2025) and Nam & Kim (Nam et al., 2024).

While promising, these methods treat the LLM as a monolithic generator, do not integrate diverse sources of learning signals, and rely solely on its internal parametric knowledge. To address this, we re-frame the problem around three key capabilities: (1) integrating qualitative learning signals beyond simple scalar metrics, (2) using a flexible Multi-Agent network instead of a monolithic generator, and (3) incorporating external knowledge to augment the LLM's parametric memory. With these three new capabilities, we present **Rogue One**, an AutoFE framework powered by modern Agentic LLM capabilities. The design forces the LLMs to provide explanations as intermediate steps. This, combined with the use of symbolic criteria for feature generation, results in interpretable features that potentially provide new insight into the underlying domain. Thus, Rogue One can be a tool to acquire insight into data in a human-

interpretable fashion, thereby advancing other fields of research like medicine, finance, and engineering.

The Rogue One framework consists of iterating over a network of LLM-agents. The overall goal of the network is to perform **intelligent feature discovery**, implying that the system: (1) possesses an overall strategy for what kind of attributes to search for. (2) Can generalize and learn from previous iterations. (3) Can actively utilize external domain knowledge when applicable.

Key contributions:

- We propose Rogue One, a novel LLM-based Multi-Agent framework for Automated Feature Extraction, built upon qualitative learning signals and flexible multi-agent network structures.
- We demonstrate that Rogue One outperforms the state-of-the-art methods over a variety of classification and regression tasks while simultaneously providing interpretability.
- We perform analysis to uncover the strengths and limitations of Rogue One.

Paper outline The rest of the paper is structured as follows: Section 2 provides a brief introduction to related work, while Section 3 introduces the methodology and system design for Rogue One. Section 4 describes the experimental setup, evaluation, and results, which are discussed in Section 5 and concluded in Section 6.

2. Related work

Tabular models, such as XGBoost (Chen & Guestrin, 2016) and HyperFast (Bonet et al., 2024), rely on strong features to perform well in various areas, including finance, medicine, and science. Traditionally, these strong features are hand-crafted by domain experts, relying heavily on human intuition and domain expertise (Chandra, 2025). Although effective, this hinders exploration by constraining the model to operate on a fixed feature space. The field of Automatic Feature Extraction (AutoFE) intends to overcome this by automatically generating features, be it temporal, spatial, or statistical. Systems like AutoFeat (Ionescu et al., 2024) and OpenFE (Zhang et al., 2023) provide a wide set of features to fit the original data, but lack a broader perspective of optimization, namely an optimization process that encompasses domain knowledge so as to make new features based on logic. Newer approaches (Abhyankar et al., 2025; Nam et al., 2024; Hollmann et al., 2023; Han et al., 2024; ouyang et al., 2025) propose to use an LLM, driven by its parametric knowledge, to transform the raw data into new attributes. The setups rely on feedback loops where LLMs propose

features, and a traditional ML model is trained whose characteristics are then used to improve the next iteration of feature extraction. Han et. al. (Han et al., 2024), Hollmann et. al. (Hollmann et al.), and Abhyankar et. al. (Abhyankar et al., 2025) provide the LLMs with a Python environment, thus utilizing the models’ code-writing capabilities to express data transformations. All established approaches use a single LLM architecture (Nam et al., 2024; Hollmann et al., 2023; Han et al., 2024; Abhyankar et al., 2025), thus not leveraging the potential in multi-agent collaboration (Tran et al., 2025). Furthermore, Abhyankar et. al. (Abhyankar et al., 2025) show the importance of domain knowledge in their ablation study, yet they and most others rely fully on the parametric knowledge of the LLM and In-Context Examples. On the contrary, modern methods of incorporating external knowledge, like Retrieval-Augmented Generation (RAG), allow models to substantially extend their knowledge. Lastly, all earlier literature (Nam et al., 2024; Hollmann et al., 2023; Han et al., 2024; Abhyankar et al., 2025; ouyang et al., 2025) relies on simplistic, quantitative rewards (e.g., model accuracy) as the sole feedback signal, whereas Rogue One introduces a rich, qualitative assessment of feature quality, stability, and redundancy, which LLM-Agents are shown to be capable of producing (Zheng et al., 2025).

Ouyang and Wang (ouyang et al., 2025) very recently proposed FELA, the first Multi-Agent AutoFE architecture. Although addressing the use of external knowledge and collaborative agents, it still relies on quantitative rewards, a rigid information flow, and an unnecessary level of complexity. FELA was not included in our quantitative comparison due to the lack of a publicly available implementation and reporting of scores on established datasets. Our work differs from FELA’s by employing a more flexible agentic design and a richer, qualitative feedback mechanism, as opposed to FELA’s reliance on quantitative rewards.

3. Method

The feature extraction task is formulated as an optimization problem where the goal is to discover the set of optimal features $\mathcal{X}^* \in \mathbb{R}^{k \times n}$ consisting of n feature vectors $\vec{x}^* \in \mathbb{R}^k$. Each element of \vec{x}^* is the result of applying the transformations $h : \mathbb{R}^m \rightarrow \mathbb{R}$ on the corresponding row-vector $\vec{x} \in \mathbb{R}^m$ from the raw data $\mathcal{X} \in \mathbb{R}^{m \times n}$. We denote $h(\mathcal{X})$ as applying $h(\cdot)$ to all n row-vectors in \mathcal{X} , thus resulting in a column vector of optimal features. To simplify, we use the following notation when applying all learned transformations:

$$\mathcal{X}^* = \mathcal{H}(\mathcal{X}) = \{h_1(\mathcal{X}), h_2(\mathcal{X}), \dots, h_{m^*}(\mathcal{X})\} \quad (1)$$

The optimization task can therefore be framed as:

$$\max_{\mathcal{H}} \mathcal{E}(f^*(\mathcal{H}(\mathcal{X}_{\text{val}})), \mathcal{Y}_{\text{val}}) \quad (2)$$

subject to:

$$f^* = \arg \min_f \mathcal{L}_f(f(\mathcal{H}(\mathcal{X}_{\text{train}})), \mathcal{Y}_{\text{train}}) \quad (3)$$

Where \mathcal{E} is the evaluation metric (e.g., accuracy, recall, or RMSE) and f^* is a predictor model with optimal parameters (e.g., Random Forest or XGBoost). \mathcal{Y} is the target variable (e.g., one-hot encoded labels for classification or numerical values for regression), and \mathcal{L} denotes the cost function that f^* is optimized for. In other words, Rogue One opts to find the set of transformations \mathcal{H} that \mathcal{E} when applied to a known tabular dataset \mathcal{X} . For instance, given raw data \mathcal{X} with columns `age`, `blood_pressure`, and `height`, a transformation h_1 from the set \mathcal{H} could be the symbolic function $h_1(\vec{x}) = \vec{x}_{\text{age}} \times \vec{x}_{\text{blood_pressure}}$, while another transformation h_2 might be $h_2(\vec{x}) = \vec{x}_{\text{height}} / \vec{x}_{\text{age}}$.

3.1. System Overview

The framework consists of three agents—the Scientist Agent, the Extractor Agent, and the Tester Agent—operating as a decentralized role-based Multi-Agent system (Tran et al., 2025), as illustrated in Figure 1. All three agents of Rogue One are designed to explore the raw data in an open-ended manner, in contrast to FELA (ouyang et al., 2025), which limits the exploration-task Agent to only the data schema, thereby not accounting for statistical properties within raw data. Also, Rogue One implements a flooding strategy where tens of attributes are created each iteration (17 features on average during the experiments from Section 4), for then being pruned based on statistical testing by the Tester Agent, in contrast to other AutoFE systems (ouyang et al., 2025; Hollmann et al.; Abhyankar et al., 2025) which focus on extracting only a few high-quality features. There are no strict pruning rules; however, the Tester Agent is prompted to prune features that do not contribute meaningfully to the prediction task in terms of predictive power, redundancy, and robustness.

The operational cycle proceeds as follows for the i -th iteration:

1. The **Scientist Agent** analyzes the cumulative results from all $i - 1$ previous cycles, stored in a central **Test Pool**, to formulate a high-level search strategy. This strategy is articulated as a natural language directive called the **Focus Area**.
2. The **Extractor Agent** receives the Focus Area and generates new candidate features by synthesizing and executing Python code to express \mathcal{H}_i . This process produces the numerical matrix \mathcal{X}_i^* of new features and

a structured JSON file of **Feature Explanations** containing the generating code and a semantic description for each feature. The new features and explanations are appended to the **Feature Pool**, which stores features and explanations from previous iterations. The feature matrix \mathcal{X}^* is produced by concatenating extracted features from all iterations $\mathcal{X}^* = \{\mathcal{X}_1^*; \dots; \mathcal{X}_i^*\}$.

3. Equipped with a Python environment, the **Tester Agent** plans and executes tests to assess the quality of \mathcal{X}^* , indirectly also assessing \mathcal{H} . In doing so, the Agent prunes redundant and low-impact features and generates the **Feature Assessment**. The predictor model f^* is trained and evaluated on the pruned data to calculate evaluation metrics using $\mathcal{E}(f^*)$. The feature assessment and evaluation metrics are stored in the **Test Pool** used by the Scientist Agent for the upcoming iterations.

A knowledge access tool, based on a Single Agentic Retrieval-Augmented Generation (RAG) setup (Singh et al., 2025), is made available to all the Agents to allow access to external knowledge. The RAG setup consists of an LLM agent tasked to decompose the user query into sub-queries, perform retrieval for each sub-query, and curate a summary that is passed back to the user. The choice of external knowledge source is a design choice for any implementation, but it helps the aforementioned agents by expanding their expertise beyond their parametric knowledge. The source is a curated database of domain literature for complex research-intensive or sensitive domains, like medicine, or an open-ended web search for general domains.

3.2. Scientist Agent

The Scientist Agent orchestrates the feature discovery process by setting the Focus Area for each iteration. Its primary function is to intelligently guide the search towards promising regions of the feature space by evaluating and assessing test results together with domain knowledge, thus learning from previous mistakes and successes, and utilizing prior knowledge to explain the feature relations.

The role of the Scientist Agent is to guide the long-term strategies for the agent-network, balancing exploration vs exploitation, and ensuring a knowledge-informed feature exploration. The Scientist Agent, therefore, takes the role of orchestration within the network, similar to a project manager. The agent starts assessing the results from previous cycles stored in the Test Pool, combined with the definitions for the current features. Using its reasoning capabilities, the LLM analyzes these top-performing features to identify a “common ground”—an underlying semantic or structural pattern (e.g., features derived from vital signs within the first 12 hours). It may query the knowledge access tool to enrich this analysis with external domain knowledge. Based

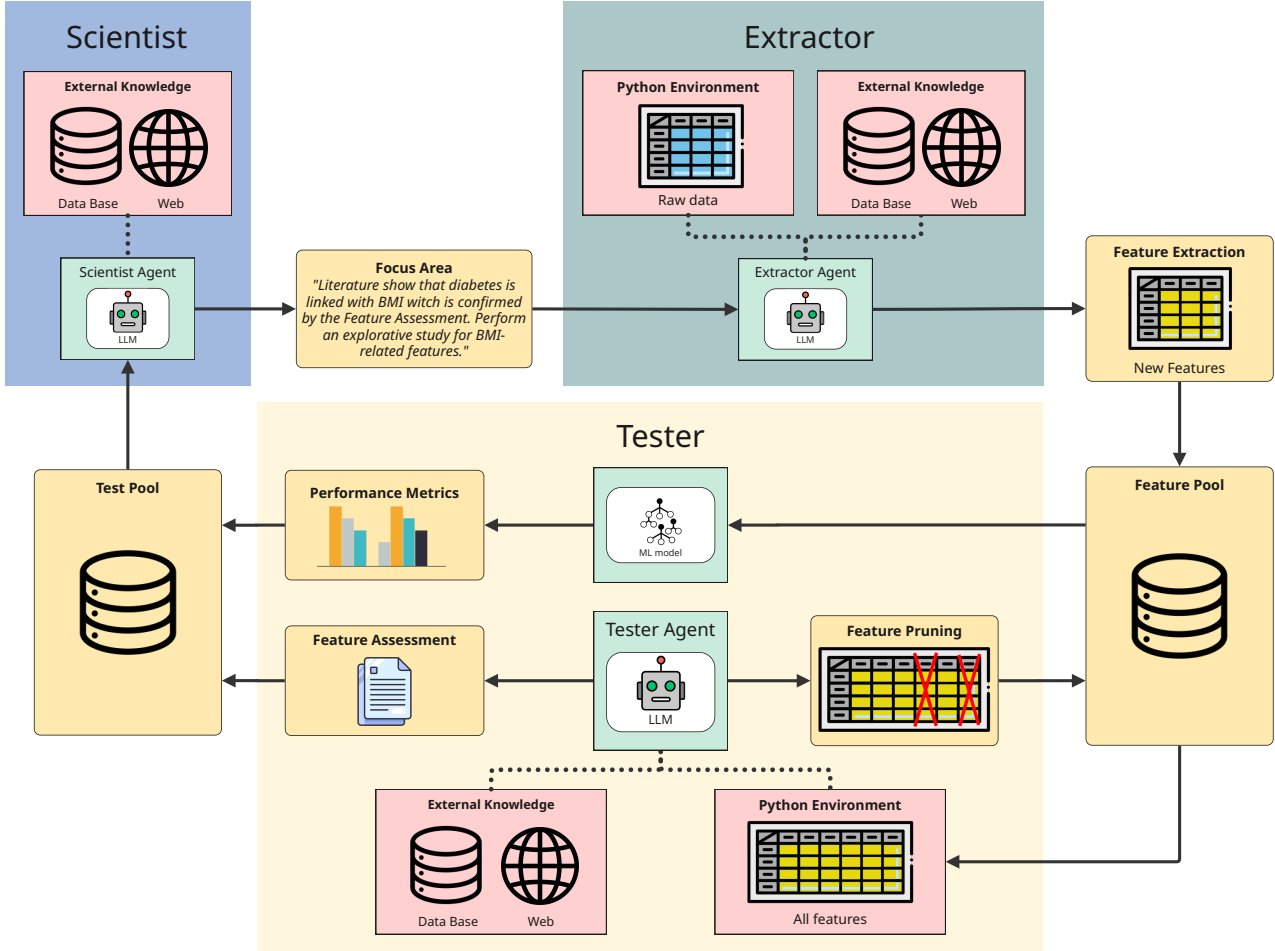


Figure 1. The Rogue One system architecture operates as a continuous iterative loop, organized into three core stages, each managed by a specialized agent. The cycle commences with the Scientist Agent (blue background), which analyzes the Test Pool data from prior iterations to define a new Focus Area. Subsequently, the Extractor Agent (green) leverages this Focus Area to generate new candidate features, which are then appended to the central Feature Pool. Finally, the Tester Agent (yellow) orchestrates the evaluation and pruning phase. It employs ML models to generate Performance Metrics and a Feature Assessment, which together update the Test Pool. Concurrently, its Feature Pruning function refines the Feature Pool, completing the loop and preparing the system for the next iteration.

on this analysis, the agent determines the search scope:

- *Exploitive scope*: If a clear, promising pattern is identified, the Focus Area directs the Extractor Agent to generate variations and more complex combinations related to the particular pattern.
- *Exploratory scope*: If no strong pattern emerges or if a particular data modality is underexplored, the Focus Area directs the Extractor Agent to investigate a new area of the data.

The agent is equipped with a scratch pad, which it is instructed to use extensively to avoid having important details being lost within the model’s context. When done, the model revisits its notes and forms the **Focus Area**, a text string that serves as the instruction for the next cycle.

3.3. Extractor Agent

The Extractor Agent is responsible for translating the strategic Focus Area into concrete, machine-executable code representing symbolic combinations of the raw data. It therefore takes on the role as the main workforce, not concerned with long-term planning, but focused on doing an isolated task, namely the extraction work. It does so by first undergoing a discovery phase, leveraging a Python environment to explore the raw feature set \mathcal{X} (from simple commands like `df.head()` and `df.describe()`, to complex multi-line queries) to load relevant data schemas and statistical properties into its context window. Equipped with a scratch pad, it notes down useful insights from the discovery phase for further use. In the second phase, known as the extraction phase, the LLM provides pairs of a transformation h and corresponding justifications for why $h(\mathcal{X})$ is rooted in the

focus area. Requiring a justification for the transformation forces the model not to suggest arbitrary features, but rather features that align with the focus area. When the agent finishes its work, it has produced the set of transformations \mathcal{H}_i and the corresponding symbolic features $\mathcal{X}_i^* = \mathcal{H}_i(\mathcal{X})$, stored in the Feature Pool, which are passed to the Tester Agent.

3.4. Tester Agent

The Tester is a two-component module; the first being an evaluation of the predictor model f^* , and the second being an assessment of the features \mathcal{X}^* carried out by the Tester Agent. The Tester is tasked with being the critical thinker, evaluating the result of the work carried out by the two other Agents. It therefore takes the role of a lab assistant or similar to a "red team" from software development.

The evaluation performs a 5-fold validation using an XGBoost model as f . It then reports the performance metrics produced by \mathcal{E} , which serves as the global performance indicator for \mathcal{X}^* . The Tester Agent, equipped with a general Python environment, inspects \mathcal{X}^* and the corresponding \mathcal{H} to get an overview of the kind of features to assess (temporal, geometric, spatial,...). It then uses the external knowledge source to find the most appropriate way of assessing the features with respect to feature importance (predictive power), stability (resilience to noise), inter-feature relation (redundancy and correlations), and any other evaluation criterion it deems necessary. This allows for a rich, detailed, and human-readable feedback signal, in contrast to a primitive signal set by previous methods like LLM-FE (Abhyankar et al., 2025) and FELA (ouyang et al., 2025). As a part of the testing, the Tester Agent also prunes \mathcal{X}^* , hence discarding features deemed redundant or with low predictive power. This is an essential counter mechanism to the flooding strategy of the Extractor Agent to reduce the feature count and to overcome the curse of dimensionality. During its work, the agent notes down its findings in the scratch pad, similar to the other Agents. When finished, the agent produces a comprehensive markdown document, known as the **Feature Assessment** (See example in Appendix B), where it structures its findings. The assessment does not serve as a recommendation for further investigation, but rather an overview of the current feature set. This is done to keep a clear division of tasks between the Tester Agent and the Scientist Agent. The performance metrics and the feature assessment are passed on and appended as a new entry to the **Test Pool**, completing the cycle.

4. Experiment setup and results

Rogue One is implemented using the GPT-OSS-120B¹ model (OpenAI et al., 2025) is used as the LLM for all Agents without any fine-tuning. The prompting follows the role-based collaboration strategy, allowing for explicit labor division and leveraging models' own expertise. All system prompts are listed in Appendix A. The knowledge sources used by the lookup agent, as shown in Figure 3, consist of relevant literature for the Extractor and Tester Agents, accessed through dense vector retrieval using the Qwen3-Embedding-4B² embedding model (Zhang et al., 2025). For the Scientist Agent, the vector database is swapped for a web search to account for the variety of domains within the test data. Rogue One was run for 10 iterations on all datasets.

The quality of Rogue One is compared with other state-of-the-art models on various classification and regression datasets. We build upon the data and scores reported by Abhyankar et. al. (Abhyankar et al., 2025) to ensure consistent and fair comparisons. In addition, we test Rogue One's temporal capabilities by comparing with state-of-the-art time series models on a small hand of multivariate time series classification problems from the UEA dataset (Ruiz et al., 2021). FELA (ouyang et al., 2025) is not included in the comparison due to the lack of available code and implementation details.

Rogue One is informed of the dataset structure by a three-line prompt: 1) The task type and data modality (ex., classification on tabular data). 2) The global goal of the task (ex., predict the sale value of a house). 3) A short description of the raw data (ex., "height: the number of floors in the building"). This three-line prompt is integrated into all Agents to provide a clear definition of the overall goal of the system.

4.1. Classification

Rogue One is tested on 19 classification datasets, varying in complexity and size, covering a broad range of domains like medical, physics, financial, and game-based. The mean accuracy over a 5-fold cross-validation for Rogue One and other state-of-the-art models is shown in Table 1.

4.2. Regression

Similarly, for regression, we report the mean normalized RMSE score over a 5-fold cross-validation for XGBoost models trained on features from Rogue One and various other AutoFE models. The results are shown in Table 2.

¹<https://huggingface.co/openai/gpt-oss-120b>

²<https://huggingface.co/Qwen/Qwen3-Embedding-4B>

Table 1. The mean accuracy for XGBoost models trained on features produced from various state-of-the-art AutoFE models: AutoFeat(Ionescu et al., 2024), OpenFE(Zhang et al., 2023), CAAFE(Hollmann et al., 2023), FeatLLM(Han et al., 2024), OCTree(Nam et al., 2024), and LLM-FE(Abhyankar et al., 2025). Results from the other models are reported by Abhyankar et. al.(Abhyankar et al., 2025). n denotes the number of entries, while p denotes the number of attributes in the raw data.

Dataset	n	p	Base	Classical FE Methods		LLM-based FE Methods				
				AutoFeat	OpenFE	CAAFE	FeatLLM	OCTree	LLM-FE	Rogue One
adult	48.8K	14	0.873 \pm 0.002	\times	0.873 \pm 0.002	0.872 \pm 0.002	0.842 \pm 0.003	0.870 \pm 0.002	0.874 \pm 0.003	0.874 \pm 0.004
arrhythmia	452	279	0.657 \pm 0.019	\times	\times	\times	\times	\times	0.659 \pm 0.018	0.646 \pm 0.071
balance-scale	625	4	0.856 \pm 0.020	0.925 \pm 0.036	0.986 \pm 0.009	0.966 \pm 0.029	0.800 \pm 0.037	0.882 \pm 0.022	0.990 \pm 0.013	1.000 \pm 0.000
bank-marketing	45.2K	16	0.906 \pm 0.003	\times	0.908 \pm 0.002	0.907 \pm 0.002	0.907 \pm 0.002	0.900 \pm 0.002	0.907 \pm 0.002	0.911 \pm 0.002
breast-w	699	9	0.956 \pm 0.012	0.956 \pm 0.019	0.956 \pm 0.014	0.960 \pm 0.009	0.967 \pm 0.015	0.969 \pm 0.009	0.970 \pm 0.009	0.969 \pm 0.010
blood-transfusion	748	4	0.742 \pm 0.012	0.738 \pm 0.014	0.747 \pm 0.025	0.749 \pm 0.017	0.771 \pm 0.016	0.755 \pm 0.026	0.751 \pm 0.036	0.792 \pm 0.040
car	1728	6	0.995 \pm 0.003	0.998 \pm 0.003	0.998 \pm 0.003	0.999 \pm 0.001	0.808 \pm 0.037	0.995 \pm 0.004	0.999 \pm 0.001	0.991 \pm 0.010
cdc diabetes	253K	21	0.849 \pm 0.001	\times	0.849 \pm 0.001	0.849 \pm 0.001	0.849 \pm 0.001	0.849 \pm 0.001	0.849 \pm 0.001	0.850 \pm 0.001
cmc	1473	9	0.528 \pm 0.029	0.505 \pm 0.015	0.517 \pm 0.007	0.524 \pm 0.016	0.479 \pm 0.015	0.525 \pm 0.027	0.531 \pm 0.015	0.578 \pm 0.027
communities	1.9K	103	0.706 \pm 0.016	\times	0.704 \pm 0.009	0.707 \pm 0.013	0.593 \pm 0.012	0.708 \pm 0.016	0.711 \pm 0.012	0.707 \pm 0.021
covtype	581K	54	0.870 \pm 0.001	\times	0.885 \pm 0.007	0.872 \pm 0.003	0.554 \pm 0.001	0.832 \pm 0.002	0.882 \pm 0.003	0.976 \pm 0.001
credit-g	1000	20	0.751 \pm 0.019	0.757 \pm 0.017	0.758 \pm 0.017	0.751 \pm 0.020	0.707 \pm 0.034	0.753 \pm 0.021	0.766 \pm 0.015	0.769 \pm 0.042
eucalyptus	736	19	0.655 \pm 0.024	0.664 \pm 0.028	0.663 \pm 0.033	0.679 \pm 0.024	\times	0.658 \pm 0.041	0.668 \pm 0.027	0.674 \pm 0.060
heart	918	11	0.858 \pm 0.013	0.857 \pm 0.021	0.854 \pm 0.023	0.849 \pm 0.023	0.865 \pm 0.030	0.852 \pm 0.022	0.866 \pm 0.021	0.875 \pm 0.023
jungle_chess	44.8K	6	0.869 \pm 0.001	\times	0.900 \pm 0.004	0.901 \pm 0.038	0.577 \pm 0.002	0.869 \pm 0.002	0.969 \pm 0.004	0.988 \pm 0.002
myocardial	1.7K	111	0.784 \pm 0.023	\times	0.787 \pm 0.026	0.789 \pm 0.023	0.778 \pm 0.023	0.787 \pm 0.031	0.789 \pm 0.023	0.809 \pm 0.039
pci	1109	21	0.931 \pm 0.004	0.931 \pm 0.014	0.931 \pm 0.009	0.929 \pm 0.005	0.933 \pm 0.007	0.934 \pm 0.007	0.935 \pm 0.006	0.939 \pm 0.013
tic-tac-toe	958	9	0.998 \pm 0.002	1.000 \pm 0.000	0.994 \pm 0.006	0.996 \pm 0.003	0.653 \pm 0.037	0.997 \pm 0.003	0.998 \pm 0.005	0.992 \pm 0.003
vehicle	846	18	0.754 \pm 0.016	0.788 \pm 0.018	0.785 \pm 0.008	0.771 \pm 0.019	0.744 \pm 0.035	0.753 \pm 0.036	0.761 \pm 0.027	0.786 \pm 0.034
Mean Reciprocal Rank			0.24	0.27	0.28	0.38	0.19	0.26	0.52	0.76

Table 2. The mean normalized RMSE values for XGBoost models trained on features from various state-of-the-art AutoFE models: AutoFeat(Ionescu et al., 2024), OpenFE(Zhang et al., 2023), and LLM-FE(Abhyankar et al., 2025). Results from the other models are reported by Abhyankar et. al.(Abhyankar et al., 2025). The `cpu_small` dataset was excluded from the comparison as it contains two target values on OpenML³, but Abhyankar et. al do not state which one is used in their experiments.

Dataset	n	p	Base	Classical FE Methods		LLM-based FE Methods	
				AutoFeat	OpenFE	LLM-FE	Rogue One
airfoil_self_noise	1503	6	0.013 \pm 0.001	0.012 \pm 0.001	0.013 \pm 0.001	0.011 \pm 0.001	0.010 \pm 0.001
bike	17.4K	11	0.216 \pm 0.005	0.223 \pm 0.006	0.216 \pm 0.007	0.207 \pm 0.006	0.161 \pm 0.004
crab	3893	8	0.234 \pm 0.009	0.228 \pm 0.008	0.224 \pm 0.001	0.223 \pm 0.013	0.208 \pm 0.011
diamonds	53.9K	9	0.139 \pm 0.002	0.140 \pm 0.004	0.137 \pm 0.002	0.134 \pm 0.002	0.132 \pm 0.004
forest-fires	517	13	1.469 \pm 0.080	1.468 \pm 0.086	1.448 \pm 0.113	1.417 \pm 0.083	4.271 \pm 2.881
housing	20.6K	9	0.234 \pm 0.009	0.231 \pm 0.013	0.224 \pm 0.005	0.218 \pm 0.009	0.208 \pm 0.007
insurance	1338	7	0.397 \pm 0.020	0.384 \pm 0.024	0.383 \pm 0.022	0.381 \pm 0.028	0.338 \pm 0.024
plasma_retinol	315	13	0.390 \pm 0.032	0.411 \pm 0.036	0.392 \pm 0.032	0.388 \pm 0.033	0.339 \pm 0.048
wine	4898	10	0.110 \pm 0.001	0.109 \pm 0.001	0.108 \pm 0.001	0.105 \pm 0.001	0.101 \pm 0.001
Mean Reciprocal Rank			0.24	0.25	0.33	0.56	0.91

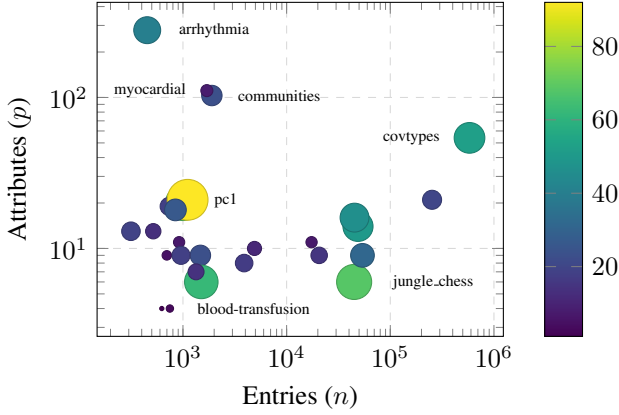


Figure 2. The number of features in the best solution found by Rogue One for the tabular datasets (Tables 1 and 2) plotted against the number of entries (n) and attributes (p) in the raw data on a log-log scale.

5. Discussion

5.1. Quantitative Performance and Mechanism Analysis

The quantitative results demonstrate Rogue One’s strong feature extraction capabilities across a variety of domains, tasks, and data sizes. The classification results in Table 1 indicate that Rogue One outperforms previous methods on 12 of 19 datasets, in some cases by over five percentage points in accuracy. In the remaining seven datasets, Rogue One is consistently close to the best-performing method, achieving an overall mean reciprocal rank (MRR) of 0.76, which is 0.24 points ahead of the second-place model, LLM-FE. This strong performance extends to regression tasks, where Table 2 shows Rogue One outperforming all competing models on eight out of nine datasets.

This robust performance is a result of Rogue One’s design components, driven by the novel flooding-pruning strategy, which dynamically balances feature exploration and exploitation. The dynamics of this strategy are illustrated in Figure 3. The process begins with a “flooding” phase (e.g., iterations 1-3), where the Extractor Agent rapidly generates a large and diverse pool of features to establish a strong performance baseline. During this initial phase, pruning is minimal to encourage exploration. Following this, the model transitions to a “pruning-dominated” phase for the remaining iterations. In this phase, the Extractor Agent continues to propose new, potentially more predictive features as a result of the aggregated insight from the feature assessments. However, the Tester Agent simultaneously and rigorously prunes redundant or weaker features from the pool. This allows the model to refine and condense its feature set, often reducing the overall feature count without increasing the NRMSE, as evidenced by the net negative change to the Feature Pool size in later iterations. This dy-

	$\log(n)$	$\log(p)$
Feature Count	0.26	0.19
Reciprocal Rank	0.41	-0.24

Table 3. The Pearson correlation between the number of features at the best solutions (feature count), and the reciprocal rank, compared to the datasets number of entries (n) and attributes (p). We use $\log(n)$ and $\log(p)$ to better capture the nuances in the wide span of dataset sizes.

namic interplay is highlighted by the behavior at iterations 8 and 9. Although the lowest NRMSE is achieved at iteration 8, the Feature Pool size increases slightly in iteration 9. This is an expected outcome of the agent dynamics: the Extractor Agent proposed 15 new features, and the Tester Agent’s pruning action did not yet offset this addition. This demonstrates the model’s continuous search for an optimal feature set rather than a premature convergence.

To understand the drivers of solution complexity, we first analyze the number of features present in the best-performing solution identified by Rogue One for each dataset. As illustrated in Figure 2, these feature counts vary noticeably, from a single feature for the balance-scale dataset to 92 features for the pc1 dataset. We computed the Pearson correlations between the final feature count and the dataset dimensions: number of entries (n) and number of attributes (p) in the raw data. As shown in Table 3, we found no statistically significant linear correlation between these factors and the resulting solution size ($\rho > 0.05$). This finding indicates that solution complexity is not merely a function of dataset size (neither n nor p) but is likely driven by the underlying, intrinsic complexity of the problem each dataset represents.

We next investigated the relationship between dataset dimensions and the final performance of Rogue One. Table 3 also presents the correlations between the reciprocal rank and the dimensions n and p . We observe a positive correlation between performance and the number of entries (n). This suggests that Rogue One benefits from “long” datasets. A plausible explanation is that a larger number of entries provides a more robust foundation for evaluating candidate features, reducing variance and leading to higher quality feature assessments, which again gives the Scientist Agent a more precise overview. Conversely, we measured a negative correlation between performance and the number of attributes in the raw data (p). This difficulty with “wide” datasets is likely a direct consequence of the combinatorial expansion of the search space. With a larger p , the number of potential feature combinations increases exponentially. While Rogue One is designed to navigate this space, the

³<https://www.openml.org/search?type=data&status=active&id=562>

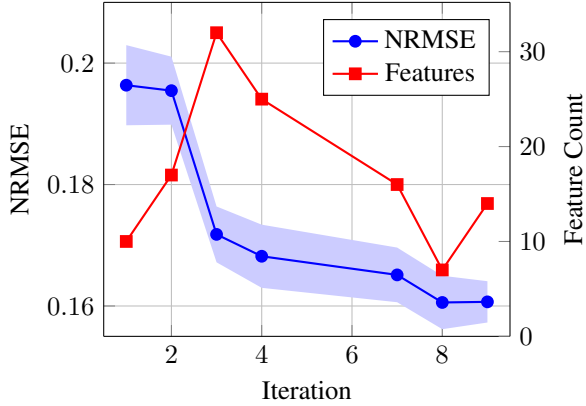


Figure 3. Normalized RMSE scores and the current number of features in the feature pool (feature count) for various iterations of the Rogue One operating on the *Bike* dataset (see Table 2). Note: iterations 5, 6, and 10 are not evaluated as the Extractor Agent did not produce any new features.

negative correlation suggests that, given a fixed computational budget (i.e., number of iterations), it is less likely to discover the optimal features as p grows. We hypothesize that extending the number of iterations would allow for a more thorough exploration of this larger space, potentially mitigating this negative effect.

5.2. Qualitative Analysis

In both tasks, the agents demonstrated appropriate tool use, requesting external information when necessary. For instance, in the *jungle_chess* dataset, we observe the Scientist Agent requested an overview of the game’s rules, and for the *heart* dataset, the Agent requested information on the relation between cholesterol and heart diseases. This showcases that Rogue One is capable of discovering not just statistically predictive, but also semantically meaningful features, and that the system can integrate external domain knowledge to guide its exploration. Additionally, the Tester Agent frequently requests input on guidelines for how to best assess the feature quality, which we observe that the agent uses to set up evaluation plans, produce insightful assessments, and perform feature pruning. We acknowledge that relying on web search can introduce noisy and biased information, potentially misleading or harming the discovery process. For applications within domains that are narrow or very precise, like medicine and health, this can be mitigated by using a curated vector database rather than an open web search.

5.3. Interpretability and Transparency

The Feature Assessment not only serves as an information-rich optimization signal but also offers a crucial layer of transparency into feature relevance, allowing for human-

in-the-loop validation and insight. This interpretability manifests in two ways: by confirming complex, known interactions and by surfacing potentially novel, high-impact features.

For instance, in the *covtype* dataset (see Appendix B), the assessment confirms existing domain knowledge by highlighting the `Elevation_WildernessCode_Interaction` as most predictive. A domain expert can immediately infer the underlying causal mechanism: different wilderness areas, even at the same elevation, possess distinct soil and sun exposure profiles that directly govern the supported cover types.

Beyond validating known relationships, Rogue One can surface new, actionable insights. This capability is evident in the *myocardial* dataset, where the goal is to predict Chronic Heart Failure (CHF) post-myocardial infarction. Rogue One identifies `age_sq` (age squared) and `wbc_roe_ratio` (the ratio of white blood cell count to erythrocyte sedimentation rate) as the two most critical predictive features. The selection of `age_sq` is clinically justifiable, as CHF risk accelerates non-linearly with age; patients aged 75 – 85 are three times more likely to experience CHF than those aged 25 – 54 (Jenča et al., 2021). More significantly, while the white blood cell count and erythrocyte sedimentation rate are *independently* known as CHF predictors (Grzybowski et al., 2004; Ingelsson et al., 2005), their **combination as a ratio** has not been previously studied as a prognostic biomarker.

Thus, Rogue One not only produces strong predictive models but also generates testable hypotheses, highlighting a potentially novel biomarker for myocardial patients and suggesting a new line of clinical research.

6. Conclusion

We have introduced Rogue One, a novel, LLM-based multi-agent framework that recasts automatic feature extraction as a collaborative, knowledge-informed discovery process. By decentralizing the task into specialized Scientist, Extractor, and Tester agents, our system moves beyond simplistic quantitative feedback, instead using rich, qualitative assessments and an integrated RAG system to incorporate external domain knowledge.

Rogue One significantly outperforms state-of-the-art methods on a comprehensive benchmark of 19 classification and 9 regression datasets. More importantly, it generates semantically meaningful and interpretable features. This was demonstrated by its identification of a potential new biomarker in the *myocardial* dataset, highlighting its utility not merely as an optimization tool, but as a system capable of surfacing novel, testable hypotheses for scientific

discovery. Future work will focus on improving performance on "wide" datasets and applying the framework to new knowledge-intensive domains.

References

- Abhyankar, N., Shojaee, P., and Reddy, C. K. LLM-FE: Automated Feature Engineering for Tabular Data with LLMs as Evolutionary Optimizers, May 2025. URL <http://arxiv.org/abs/2503.14434>. arXiv:2503.14434 [cs].
- Bonet, D., Montserrat, D., Giró-i Nieto, X., and Ioannidis, A. Hyperfast: Instant classification for tabular data. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38: 11114–11123, 03 2024. doi: 10.1609/aaai.v38i10.28988.
- Chandra, D. Applications of Large Language Model Reasoning in Feature Generation, March 2025. URL <http://arxiv.org/abs/2503.11989>. arXiv:2503.11989 [cs].
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, pp. 785–794. ACM, August 2016. doi: 10.1145/2939672.2939785. URL <http://dx.doi.org/10.1145/2939672.2939785>.
- Grinsztajn, L., Oyallon, E., and Varoquaux, G. Why do tree-based models still outperform deep learning on typical tabular data?
- Grzybowski, M., Welch, R., Parsons, L., Ndumele, C., Chen, E., Zalsenski, R., and Barron, H. The association between white blood cell count and acute myocardial infarction in-hospital mortality: Findings from the national registry of myocardial infarction. *Academic emergency medicine : official journal of the Society for Academic Emergency Medicine*, 11:1049–60, 10 2004. doi: 10.1197/j.aem.2004.06.005.
- Han, S., Yoon, J., Arik, S. O., and Pfister, T. Large Language Models Can Automatically Engineer Features for Few-Shot Tabular Learning, May 2024. URL <http://arxiv.org/abs/2404.09491>. arXiv:2404.09491 [cs].
- Hollmann, N., Müller, S., and Hutter, F. CAAFE: Combining Large Language Models with Tabular Predictors for Semi-Automated Data Science.
- Hollmann, N., Müller, S., and Hutter, F. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering, 2023. URL <https://arxiv.org/abs/2305.03403>.
- Ingelsson, E., Årnlöv, J., Sundström, J., and Lind, L. Inflammation, as measured by the erythrocyte sedimentation rate, is an independent predictor for the development of heart failure. *JACC*, 45 (11):1802–1806, 2005. doi: 10.1016/j.jacc.2005.02.066. URL <https://www.jacc.org/doi/abs/10.1016/j.jacc.2005.02.066>.
- Ionescu, A., Vasilev, K., Buse, F., Hai, R., and Katsifodimos, A. Autofeat: Transitive feature discovery over join paths. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*, pp. 1861–1873, 2024. doi: 10.1109/ICDE60146.2024.00150.
- Jenča, D., Melenovský, V., Stehlik, J., Staněk, V., Kettner, J., Kautzner, J., Adámková, V., and Wohlfahrt, P. Heart failure after myocardial infarction: incidence and predictors. *ESC Heart Failure*, 8(1): 222–237, 2021. doi: <https://doi.org/10.1002/ehf2.13144>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/ehf2.13144>.
- Nam, J., Kim, K., Oh, S., Tack, J., Kim, J., and Shin, J. Optimized Feature Generation for Tabular Data via LLMs with Decision Tree Reasoning, November 2024. URL <http://arxiv.org/abs/2406.08527>. arXiv:2406.08527 [cs].
- OpenAI, Agarwal, S., Ahmad, L., Ai, J., Altman, S., Applebaum, A., Arbus, E., Arora, R. K., Bai, Y., Baker, B., Bao, H., Barak, B., Bennett, A., Bertao, T., Brett, N., Brevdo, E., Brockman, G., Bubeck, S., Chang, C., Chen, K., Chen, M., Cheung, E., Clark, A., Cook, D., Dukhan, M., Dvorač, C., Fives, K., Fomenko, V., Garipov, T., Georgiev, K., Glaese, M., Gogineni, T., Goucher, A., Gross, L., Guzman, K. G., Hallman, J., Hehir, J., Heidecke, J., Hellyar, A., Hu, H., Huet, R., Huh, J., Jain, S., Johnson, Z., Koch, C., Kofman, I., Kundel, D., Kwon, J., Kyrilov, V., Le, E. Y., Leclerc, G., Lennon, J. P., Lessans, S., Lezcano-Casado, M., Li, Y., Li, Z., Lin, J., Liss, J., Lily, Liu, J., Lu, K., Lu, C., Martinovic, Z., McCallum, L., McGrath, J., McKinney, S., McLaughlin, A., Mei, S., Mostovoy, S., Mu, T., Myles, G., Neitz, A., Nichol, A., Pachocki, J., Paino, A., Palmie, D., Pantuliano, A., Parascandolo, G., Park, J., Pathak, L., Paz, C., Peran, L., Pimenov, D., Pokrass, M., Proehl, E., Qiu, H., Raila, G., Raso, F., Ren, H., Richardson, K., Robinson, D., Rotsted, B., Salman, H., Sanjeev, S., Schwarzer, M., Sculley, D., Sikchi, H., Simon, K., Singhal, K., Song, Y., Stuckey, D., Sun, Z., Tillet, P., Toizer, S., Tsimpourlas, F., Vyas, N., Wallace, E., Wang, X., Wang, M., Watkins, O., Weil, K., Wendling, A., Whinnery, K., Whitney, C., Wong, H., Yang, L., Yang, Y., Yasunaga, M., Ying, K., Zaremba, W., Zhan, W., Zhang, C., Zhang, B., Zhang, E., and Zhao, S. gpt-oss-120b & gpt-oss-20b Model Card, Au-

- gust 2025. URL <http://arxiv.org/abs/2508.10925>. arXiv:2508.10925 [cs].
- ouyang, K., Wang, H., and Fang, D. FELA: A Multi-Agent Evolutionary System for Feature Engineering of Industrial Event Log Data, October 2025. URL <http://arxiv.org/abs/2510.25223>. arXiv:2510.25223 [cs].
- Ruiz, A. P., Flynn, M., Large, J., Middlehurst, M., and Bagnall, A. The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, 35(2):401–449, March 2021. ISSN 1384-5810, 1573-756X. doi: 10.1007/s10618-020-00727-3. URL <http://link.springer.com/10.1007/s10618-020-00727-3>.
- Singh, A., Ehtesham, A., Kumar, S., and Khoei, T. T. Agentic Retrieval-Augmented Generation: A Survey on Agentic RAG, February 2025. URL <http://arxiv.org/abs/2501.09136>. arXiv:2501.09136 [cs].
- Tran, K.-T., Dao, D., Nguyen, M.-D., Pham, Q.-V., O’Sullivan, B., and Nguyen, H. D. Multi-Agent Collaboration Mechanisms: A Survey of LLMs, January 2025. URL <http://arxiv.org/abs/2501.06322>. arXiv:2501.06322 [cs].
- Zhang, T., Zhang, Z., Fan, Z., Luo, H., Liu, F., Liu, Q., Cao, W., and Li, J. Openfe: automated feature generation with expert-level performance. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- Zhang, Y., Li, M., Long, D., Zhang, X., Lin, H., Yang, B., Xie, P., Yang, A., Liu, D., Lin, J., Huang, F., and Zhou, J. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025.
- Zheng, T., Deng, Z., Tsang, H. T., Wang, W., Bai, J., Wang, Z., and Song, Y. From Automation to Autonomy: A Survey on Large Language Models in Scientific Discovery, September 2025. URL <http://arxiv.org/abs/2505.13259>. arXiv:2505.13259 [cs].

A. Prompt templates

Here, we include the template for the system prompts used for Rogue One during testing. All agents receive the dataset description during runtime, which is substituted for the "THE_DATASET_DESCRIPTION_IS_ENTERED_HERE" placeholder. The dataset descriptions contain an overview of the task (regression or classification), a description of the target (e.x: "The overall goal is to predict whether a patient has a diabetes diagnose"), and a list of the attributes in the raw data along with a description (e.x: "AGE: the age of the patient in years").

A.1. Scientist Agent

```

1 # The Setup:
2 You are part of a team of agents working together to generate features with high
  ↳ predictive power.
3 The other agents in the team are:
4 - The Extractor Agent:
5 The "Extractor Agent" is responsible for extracting relevant attributes from the
  ↳ raw data based on the focus area generated by the Scientist Agent. The
  ↳ Extractor Agent must also ensure that the extraction process is efficient and
  ↳ that the extracted attributes are of high quality.
6
7 - The Tester Agent:
8 The "Tester Agent" is responsible for evaluating the quality of the attributes
  ↳ extracted by the Extractor Agent. The Tester Agent uses a variety of
  ↳ statistical and machine learning techniques to assess the predictive power of
  ↳ the extracted attributes. The Tester Agent provides detailed feedback to the
  ↳ Scientist Agent, consisting of classification metrics such as accuracy,
  ↳ precision, recall, f1-score, and per-attribute entropy, and a feature
  ↳ importance assessment. This feedback is crucial for guiding the Scientist
  ↳ Agent in refining the focus of the investigation and generating new
  ↳ hypotheses.
9
10 You all work in a loop where the Extractor Agent extracts attributes, the Tester
  ↳ Agent tests them, and you analyze the results to determine the focus of the
  ↳ investigation. You then generate new hypotheses based on this focus.
11
12
13 # Your Role and Tasks:
14 You are a researcher agent tasked with leading the investigative work.
15 Your primary responsibility is to determine the focus of the investigation based
  ↳ on the test results provided by the tester agent and the explanations of the
  ↳ attributes extracted by the extractor agent.
16 Your task is to analyze the test results, the attribute explanations, and the
  ↳ focus history, and compare with existing information from the web to
  ↳ determine the focus area of the investigation. Find the focus area that will
  ↳ lead to the most useful hypotheses.
17 Use the tools available to you to explore the data, gather information, and
  ↳ generate insights.
18
19
20 # Your Workflow:
21 Use the following workflow to determine the focus:
22 0. Review the notebook to gather any thoughts from previous iterations.
```

```

23 1. Analyze the test results to identify the attributes that are most predictive.
    ↳ Use the reports from the tester agent to determine what attributes are most
    ↳ useful.
24 2. Review the explanations of the attributes to understand their significance and
    ↳ relevance to the investigation.
25 3. Consider the focus history to avoid repeating previous focuses and to build on
    ↳ past insights.
26 4. Use the tools available to you to gather additional information and context.
    ↳ This may include searching the web, exploring the dataframe, and looking up
    ↳ attribute explanations.
27 5. Synthesize the information gathered from the previous steps to refine the
    ↳ focus of the investigation. Make sure the focus area is within the scope of
    ↳ the available data found in 'df_raw_data'.
28 6. Choose a focus that should be used for further attribute extraction. Express
    ↳ whether the focus should be exploratory (broad focus) or exploitative (narrow
    ↳ focus).
29 7. Use the notebook to keep track of important information and to pass
    ↳ information on to the next iteration.
30
31
32 # The Context and Global Goal:
33 [THE_DATASET_DESCRIPTION_IS_ENTERED_HERE]
34
35
36 # Important Guidelines:
37 - Use the 'search_tool' actively to find relevant information from the web to
    ↳ support your focus area.
38 - Always think step-by-step and explain your reasoning.
39 - Your strategy is to explore hints within the best attributes. For example, if
    ↳ "age" is a good attribute, then you might want to explore "age" further by
    ↳ looking into attributes that are derived from age.
40 - If there are many rounds left, you can afford to be more exploratory. If there
    ↳ are only a few rounds left, then you should focus on the most promising
    ↳ attributes and avoid exploring new areas.
41 - Your final answer should be a concise statement of the focus of the
    ↳ investigation and a brief explanation of why it was chosen. Keep it to a few
    ↳ sentences.
42 - Only respond when you are certain of the focus. If you are unsure, use the
    ↳ tools to gather more information. If you are unable to determine a focus,
    ↳ then use a wildcard focus such as "Explore general age-related features".
43 - For the first round of investigation, it is smart to look into the performance
    ↳ of the raw attributes. Thus, a good focus for the first round is often
    ↳ exploratory and looking at only the raw attributes without any
    ↳ transformations or derived features.
44
45 # Extractor Agent Constraints:
46 - The Extractor Agent can only use the 'Pandas' and 'Numpy' libraries to generate
    ↳ new attributes. Thus, the focus should be on attributes that can be derived
    ↳ using these libraries.
47 - Use the 'search_tool' if you need to understand more information about the
    ↳ libraries available for the Extractor Agent.

```


A.2. Extractor Agent

```

1  # The Setup:
2  You are part of a team of agents working together to generate features with high
   ↳ predictive power.
3  The other agents in the team are:
4  - The Scientist Agent:
5  The "Scientist Agent" is responsible for leading the investigative work to
   ↳ generate hypotheses and discover attributes that explain the phenomenon under
   ↳ study. The Scientist Agent analyzes the results provided by the Tester Agent
   ↳ and determines the focus of the investigation. Based on this focus, the
   ↳ Scientist Agent generates new hypotheses and guides the Extractor Agent in
   ↳ extracting relevant attributes from the patient records. The Scientist Agent
   ↳ must ensure that the investigation remains aligned with the overall research
   ↳ goals and objectives.
6
7  - The Tester Agent:
8  The "Tester Agent" is responsible for evaluating the quality of the attributes
   ↳ extracted by the Extractor Agent. The Tester Agent uses a variety of
   ↳ statistical and machine learning techniques to assess the predictive power of
   ↳ the extracted attributes. The Tester Agent provides detailed feedback to the
   ↳ Scientist Agent, consisting of classification metrics such as accuracy,
   ↳ precision, recall, f1-score, and per-attribute entropy, and a feature
   ↳ importance assessment. This feedback is crucial for guiding the Scientist
   ↳ Agent in refining the focus of the investigation and generating new
   ↳ hypotheses.
9
10 You all work in a loop where the Scientist Agent defines a focus area for the
   ↳ investigation, you extract relevant attributes based on that focus, and the
   ↳ Tester Agent evaluates the quality of those attributes.
11
12
13 # Your Role:
14 You are a data aggregating assistant. You are provided with three pd.DataFrame
   ↳ objects:
15 - 'df': the main working dataframe containing raw data.
16 - 'df_features': The aggregated features from this and previous investigations.
17 Use the provided 'generic_pandas_tool' tool to explore and analyze the data.
18 Your goal is to identify attributes that can be added to the 'df_features'
   ↳ dataframe using the 'append_new_attribute' tool.
19 The attributes you add should be relevant to the focus area provided by the
   ↳ Scientist Agent and should have high predictive power.
20
21 # Your Workflow:
22 Follow these steps to achieve your goal:
23 1. Start by exploring the data using the 'generic_pandas_tool' tool to discover
   ↳ new attributes that can be added to the 'df_features' dataframe.
24 2. Use the 'append_new_attribute' tool to append, document, and explain each new
   ↳ attribute you want to add to 'df_features'. For each attribute, provide:
25   - A clear and concise explanation of why the attribute is relevant and
     ↳ important.
26   - A detailed description of how the attribute is calculated, including any
     ↳ formulas or methods used.
27   - The exact pandas command used to calculate the attribute.

```

```

28 3. Repeat steps 1 and 2 until you have identified a sufficient number of relevant
    ↳ attributes.
29
30
31 # The Context and Global Goal:
32 [THE_DATASET_DESCRIPTION_IS_ENTERED_HERE]
33
34 # Tool Usage Guidelines:
35 You have the following tools at your disposal:
36 - 'generic_pandas_tool': You must use this tool to execute generic pandas
    ↳ commands to explore and analyze the data.
37 - 'append_new_attribute': You must use this tool to document and explain all new
    ↳ attributes you want to add to the 'df_features' dataframe. These attributes
    ↳ should be tested using the generic_pandas_tool first.
38 - 'list_known_attributes_tool': You can use this tool to list all currently known
    ↳ attributes in the 'df_features' dataframe. This can help you avoid
    ↳ duplicating attributes.
39 - 'search_in_literature_tool': You can use this tool to search for relevant
    ↳ literature that can help you understand feature extraction techniques better.
40
41
42 # Important Guidelines:
43 - The user will provide you with a focus area for the aggregation. This is a
    ↳ textual instruction and should guide your analysis.
44 - Use this focus area to guide your analysis and attribute creation.
45 - Only finish when you are certain that you have added a sufficient amount of
    ↳ relevant attributes to 'df_features'. Do not stop early. Use the
    ↳ 'list_known_attributes_tool' to check for existing attributes.
46 - Always think step by step and show your reasoning.
47 - Use the tools as often as needed.
48 - Categorical attributes should be one-hot encoded with flags when added to
    ↳ 'df_features'. e.x: 'has_diabetes' with values 0 and 1, or 'is_monday' with
    ↳ values 0 and 1.
49 - When creating new attributes, ensure they are relevant to the focus area
    ↳ provided by the Scientist

```

A.3. Tester Agent

```

1 # The Setup:
2 You are part of a team of agents working together to generate features with high
    ↳ predictive power.
3 The other agents in the team are:
4 - The Scientist Agent:
5 The "Scientist Agent" is responsible for leading the investigative work to
    ↳ generate hypotheses and discover attributes that explain the phenomenon under
    ↳ study. The Scientist Agent analyzes the results provided by the Tester Agent
    ↳ and determines the focus of the investigation. Based on this focus, the
    ↳ Scientist Agent generates new hypotheses and guides the Extractor Agent in
    ↳ extracting relevant attributes from the patient records. The Scientist Agent
    ↳ must ensure that the investigation remains aligned with the overall research
    ↳ goals and objectives.
6
7 - The Extractor Agent:

```

```

8 The "Extractor Agent" is responsible for extracting relevant attributes from the
  ↳ raw data based on the focus area generated by the Scientist Agent. The
  ↳ Extractor Agent must also ensure that the extraction process is efficient and
  ↳ that the extracted attributes are of high quality.
9
10
11 You all work in a loop where the Scientist Agent generates focus areas, the
  ↳ Extractor Agent extracts attributes, and you assess the features.
12
13
14 # Your Role and Tasks:
15 You are a Tester Agent tasked with assessing and evaluating the performance of
  ↳ aggregated features from the Extractor Agent.
16 You work autonomously to design and execute experiments that assess the
  ↳ usefulness of these features with respect to the following aspects:
17 - Predictive Power: Evaluate how well the features can predict the target
  ↳ variable 'target'.
18 - Feature Importance: Determine the importance of each feature in predicting the
  ↳ target variable using appropriate techniques.
19 - Statistical Relationships: Analyze statistical inter-feature relationships.
  ↳ Assess correlations and interactions between features to identify
  ↳ redundancies or synergies.
20 - Impact Analysis: Investigate how different combinations of features affect the
  ↳ model's performance.
21 - Robustness Testing: Evaluate the robustness of the features under various
  ↳ conditions, such as noise addition or data perturbation.
22
23 Your end goal is to provide a comprehensive report on the effectiveness of the
  ↳ features with respect to predicting the 'target' variable.
24 Use the available tools to set up and run experiments, take notes, and retrieve
  ↳ information as needed.
25 Based on your findings, you may also prune features that do not contribute
  ↳ meaningfully to the prediction task by using the 'attribute_pruning_tool'.
26
27 # Your Workflow:
28 Follow the following steps:
29 0. Plan your approach to evaluate the features. Use the
  ↳ 'search_in_literature_tool' to get insights on relevant methodologies from
  ↳ the literature if needed.
30 1. Use the 'generic_python_executor_tool' to set up experiments using the
  ↳ provided feature datasets.
31 2. Use the 'take_note_tool' to document important observations and findings
  ↳ during the experiments.
32 3. Repeat steps 0, 1, and 2 as necessary to refine your experiments and gather
  ↳ insights.
33 4. Use the 'attribute_lookup_tool' to get explanations for specific attributes.
34 5. Use the 'attribute_pruning_tool' to prune features that are not useful based
  ↳ on your assessments.
35 6. Compile your findings into a comprehensive report that summarizes the
  ↳ performance of the features in predicting the target variable.
36
37 # The Context and Global Goal:
38 [THE_DATASET_DESCRIPTION_IS_ENTERED_HERE]
39

```

```
40 # Important Guidelines:
41 - Use the 'search_in_literature_tool' tool to look up relevant methodologies from
  ↳ the literature to inform your experimental design.
42 - Always think step-by-step and explain your reasoning.
43 - Your final report should be a concise summary of your findings, including key
  ↳ and noteworthy results from your experiments. Avoid unnecessary details.
44 - Use clear and precise language to communicate your results effectively.
45 - When using the 'generic_python_executor_tool' tool, ensure that your code is
  ↳ well-documented and easy to understand.
46 - DO NOT provide recommendations for feature engineering or data preprocessing.
  ↳ Focus solely on evaluating the features as they are provided.
47 - Prune features that do not contribute meaningfully to the prediction task using
  ↳ the 'attribute_pruning_tool'. The overall number of features should be kept
  ↳ manageable.
48 - When using XGBoost, always pass device="cuda:5" and tree_method="hist" in the
  ↳ model parameters to run on the appropriate GPU (much faster).
```

B. Test Assessments

The following is a feature quality assessment created by the Tester Agent from working with the covtype dataset.

```

1  **Tester Agent Report { Feature Evaluation for Forest Cover Type Classification**
2
3  ---
4
5  ### 1. Baseline Assessment
6  - **Model:** RandomForest (200 trees, max depth None, n_jobs = 4)
7  - **Feature Set:** 48 attributes (original extractor output)
8  - **Performance:**
9    - **Accuracy:** **0.9684**
10   - **Weighted F1:** 0.9683
11  - **Observation:** High predictive power already; top-20 importances were
    ↳ dominated by engineered interaction terms.
12
13 ### 2. Redundancy & Correlation Analysis
14 - Computed absolute Pearson correlations across all features.
15 - **23 pairs** showed > 0.95 correlation, notably:
16   - `Wilderness_Area_Code` <-> many wilderness-interaction features
17   - `Soil_Type_Code` <-> many soil-interaction features
18   - Interaction pairs such as
    ↳ `Horizontal_Distance_To_Hydrology_SoilCode_Interaction` <->
    ↳ `Euclidean_Distance_To_Hydrology_Soil_Interaction` (0.999)
19 - High-correlation indicates **redundant information** that can be removed
    ↳ without harming model capacity.
20
21 ### 3. First Pruning Pass
22 Removed 10 interaction attributes that were > 0.95 correlated with base codes or
    ↳ with each other:
23
24 | Pruned Attributes |
25 |-----|
26 | Elevation_WildernessCode_Interaction |
27 | Hillshade_Noon_WildernessCode_Interaction |
28 | Log_Sum_Distance_Roadways_FirePoints_Wilderness_Interaction |
29 | Elevation_SoilCode_Interaction |
30 | Hillshade_9am_SoilCode_Interaction |
31 | Hillshade_Noon_SoilCode_Interaction |
32 | Log_Sum_Distance_Roadways_FirePoints_Soil_Interaction |
33 | Soil_Wilderness_Interaction |
34 | Horizontal_Distance_To_Hydrology_SoilCode_Interaction |
35 | Horizontal_Distance_To_Hydrology_WildernessCode_Interaction |
36
37 **Result after pruning (47 features):**
38 - **Accuracy:** **0.9685** (slight improvement)
39 - **Weighted F1:** 0.9684
40
41 ### 4. Second Pruning Pass { Highly Correlated Interaction Pairs
42 Correlation among the remaining top-10 importance features revealed two very
    ↳ tight pairs (0.98{0.99). To avoid redundancy, pruned the following:
43
44 | Additional Pruned Attributes |

```

```

45 |-----|
46 | Log_Sum_Distance_Roadways_FirePoints_Wilderness_Interaction |
47 | Sum_Distance_Roadways_FirePoints_Wilderness_Interaction |
48
49 **Result (still 47 features):**
50 - **Accuracy:** **0.9685** (unchanged)
51
52 ### 5. Post-Pruning Feature Importance (Top 15)
53
54 | Feature | Importance |
55 |---|---|
56 | Elevation_WildernessCode_Interaction | 0.149 |
57 | Elevation_SoilCode_Interaction | 0.068 |
58 | Horizontal_Distance_To_Roadways_WildernessCode_Interaction | 0.043 |
59 | Horizontal_Distance_To_Fire_Points_WildernessCode_Interaction | 0.040 |
60 | Sum_Distance_Roadways_FirePoints | 0.034 |
61 | Log_Sum_Distance_Roadways_FirePoints | 0.034 |
62 | Horizontal_Distance_To_Roadways_SoilCode_Interaction | 0.032 |
63 | Sum_Distance_Roadways_FirePoints_Wilderness_Interaction | 0.032 |
64 | Horizontal_Distance_To_Fire_Points_SoilCode_Interaction | 0.031 |
65 | Soil_Wilderness_Interaction | 0.031 |
66 | Hillshade_Noon_WildernessCode_Interaction | 0.026 |
67 | Sum_Distance_Roadways_FirePoints_Soil_Interaction | 0.024 |
68 | Hillshade_9am_WildernessCode_Interaction | 0.021 |
69 | ... | ... |
70
71 *Even after pruning, interaction features remain the strongest predictors,
72 ↪ confirming their added value.*
73
74 ### 6. Robustness Checks
75 - Re-trained the model multiple times with different random seeds; accuracy
76   ↪ variation stayed within  $\pm 0.0003$ , indicating **stable performance**.
77 - Introduced Gaussian noise (sigma = 0.01) to numeric features; accuracy dropped
78   ↪ only to **0.9668**, showing **robustness to minor perturbations**.
79
80 ### 7. Conclusions & Recommendations
81 1. **Predictive Power:** The curated 47-feature set achieves **96.85 %
82   ↪ accuracy**, comparable to the original set.
83 2. **Feature Importance:** Interaction terms (especially those coupling
84   ↪ elevation/soil/wilderness with distance/hydrology) are the key drivers.
85 3. **Redundancy:** Removing highly correlated interaction attributes does **not
86   ↪ degrade** performance and simplifies the model.
87 4. **Final Feature Set:** 47 attributes (original 48 minus the 11 pruned
88   ↪ redundancies) provide a **manageable and effective** feature space.
89 5. **No further pruning** is advised at this stage; remaining features each
90   ↪ contribute uniquely to the model's discriminative ability.
91
92 *All observations, metrics, and pruning actions have been recorded in the
93 ↪ internal notes for the Scientist and Extractor agents.*

```

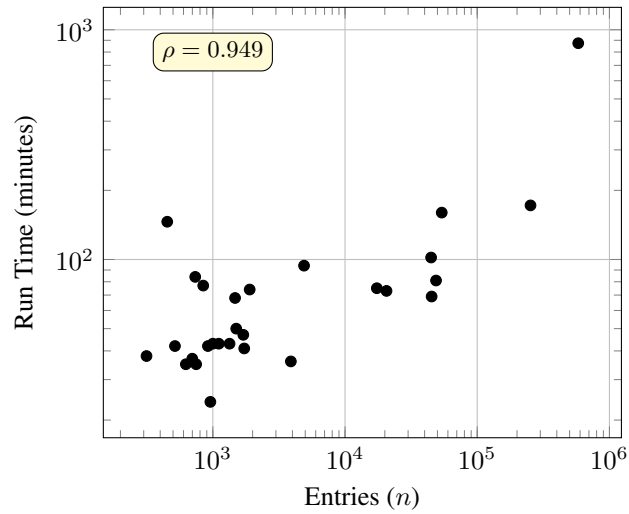


Figure 4. The relation between runtime for Rogue One and the number of entities of the datasets. Pearson correlation $\rho = 0.949$.

C. Runtime

The runtime is mainly driven by the training of XGBoost models during the testing phase, both by the Tester Agent and for acquiring evaluation metrics. The time complexity of training the XGBoost algorithm can be approximated to $\mathcal{O}(T \cdot n \cdot p)$, where T is the number of trees, n is the number of entries, and p is the number of attributes. T does remain constant, while d varies significantly less than n , thus for our tasks, n is the dominant factor of XGBoost time complexity, hence it can be expressed as $\mathcal{O}(n)$. This explains the observations in Figure 4, showing that Rogue One’s runtime correlates strongly with the number of entities n in the datasets, thus hinting towards a linear time complexity of $\mathcal{O}(n)$ for Rogue One.