AUTOKAGGLE: A MULTI-AGENT FRAMEWORK FOR AUTONOMOUS DATA SCIENCE COMPETITIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Data science tasks involving tabular data present complex challenges that require sophisticated problem-solving approaches. We propose AutoKaggle, a powerful and user-centric framework that assists data scientists in completing daily data pipelines through a collaborative multi-agent system. AutoKaggle implements an iterative development process that combines code execution, debugging, and comprehensive unit testing to ensure code correctness and logic consistency. The framework offers highly customizable workflows, allowing users to intervene at each phase, thus integrating automated intelligence with human expertise. Our universal data science toolkit, comprising validated functions for data cleaning, feature engineering, and modeling, forms the foundation of this solution, enhancing productivity by streamlining common tasks. We selected 8 Kaggle competitions to simulate data processing workflows in real-world application scenarios. Evaluation results demonstrate that AutoKaggle achieves a validation submission rate of 0.85 and a comprehensive score of 0.82 in typical data science pipelines, fully proving its effectiveness and practicality in handling complex data science tasks.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

In recent years, with the rapid development of large language models (LLMs) (OpenAI, 2022; 2023), automated data science has gradually become possible. LLM-based agents have shown great potential in the data domain, as they can automatically understand, analyze, and process data (Hassan et al., 2023; Lucas, 2023; Zhang et al., 2024a), thereby promoting the democratization and widespread application of data science.

034 However, existing research still has significant shortcomings in addressing complex data science 035 problems. Many studies are limited to simple, one-step data analysis tasks (Zhang et al., 2024c; 036 Hu et al., 2024), which are far from the actual application scenarios of data science. While recent 037 work (Jing et al., 2024) attempts to evaluate data science capabilities through more comprehensive tasks, it still focuses on relatively constrained scenarios that represent only portions of a complete data science pipeline. Other research relies on pre-built knowledge bases (Guo et al., 2024), raising 040 the barrier to use and limiting the flexibility and adaptability of solutions. Moreover, current research 041 focuses excessively on improving task completion rates and optimizing performance metrics, while 042 neglecting the interpretability and transparency of intermediate decision-making steps in logically complex data science tasks. This neglect not only affects users' understanding of solutions but also 043 diminishes their credibility and practicality in real-world applications. 044

To address these issues, we propose AutoKaggle, a universal multi-agent framework that provides
 data scientists with end-to-end processing solutions for tabular data, helping them efficiently complete daily data pipelines and enhance productivity. AutoKaggle has the following features:

(i) Phase-based Workflow and Multi-agent Collaboration. AutoKaggle employs a phase-based workflow and multi-agent collaboration system. It divides the data science competition process into six key phases: background understanding, preliminary exploratory data analysis, data cleaning (DC), in-depth exploratory data analysis, feature engineering (FE), and model-building, -validation, and -prediction (MBVP). To execute these phases, five specialized agents (Reader, Planner,

¹All code and data are available: https://anonymous.4open.science/r/AutoKaggle-B8D2.

Developer, Reviewer, and Summarizer) work collaboratively to execute these phases, from problem analysis to report generation.

 (ii) Iterative Debugging and Unit Testing. AutoKaggle ensures code quality through iterative debugging and unit testing. The Developer employs three main tools (code execution, debugging, and unit testing) to verify both syntactic correctness and logical consistency.

(iii) Machine Learning Tools Library. AutoKaggle integrates a comprehensive machine learn ing tools library covering data cleaning, feature engineering, and model-building, -validation, and
 -prediction. The library includes expert-written code snippets and custom tools, enhancing code
 generation efficiency and quality. By combining predefined tools with self-generated code, AutoK aggle handles complex tasks while reducing reliance on LLMs for domain-specific knowledge.

(iv) Comprehensive Reporting. AutoKaggle generates detailed reports after each phase and at the
 competition's conclusion, showcasing its decision-making process, key findings, actions, and reasoning. This feature makes the data processing workflows transparent, increasing user trust in AutoKaggle.

AutoKaggle provides a universal and comprehensive solution for a wide variety of data science tasks.
 By simply providing a task overview, it can automatically complete the entire process from development to testing, making it exceptionally easy to use. AutoKaggle is highly adaptable, allowing users to customize it according to their specific needs. Moreover, it offers clear interpretability throughout the automated data science process, enhancing users' understanding and trust in the system.

074 We chose competitions from the Kaggle platform to evaluate our framework. Kaggle data science 075 competitions simulate the real challenges faced by data scientists, covering the complete process 076 from data cleaning to model deployment. These competitions require participants to execute a series 077 of complex and interdependent tasks. These include: data cleaning and preprocessing, exploratory data analysis, feature engineering, and modeling. Each step demands professional knowledge and meticulous planning, often necessitating multiple iterations. This complexity makes Kaggle an ideal 079 platform for assessing the effectiveness of data science automation tools. In the 8 Kaggle data science competitions we evaluated, AutoKaggle achieved 0.85 in valid submission rate and 0.82 in 081 comprehensive score. We summarize our contributions as follows: 082

083 084

085

090

092

093

094

• We propose AutoKaggle, a novel multi-agent framework for Kaggle data science competitions, achieving high task completion rates and competitive performance above the average human level in our evaluations.

- We introduce a phase-based workflow integrated with multi-agent collaboration, incorporating iterative debugging and unit testing, which systematically addresses the complexities of data science tasks and ensures robust, correct code generation.
- We develop a machine learning tools library and integrate it into our framework, enhancing code generation efficiency and quality for complex data science tasks.
- We implement a comprehensive reporting system that provides detailed insights into the decision-making process at each phase, making AutoKaggle both a solution provider and an educational tool for data science competitions, thereby contributing to the democratization of data science skills.
- 095 096
- 097 098

099

2 AUTOKAGGLE

2.1 OVERALL FRAMEWORK

In this section, we introduce AutoKaggle, a fully automated, robust, and user-friendly framework
designed to produce directly submittable prediction results using only the original Kaggle data.
Given the diversity of data science problems, the range of potential solutions, and the need for
precise reasoning and real-time understanding of data changes, effectively handling complex data
science tasks on Kaggle is challenging. Our technical design addresses two primary issues: (*i*) how
to decompose and systematically manage complex data science tasks; and (*ii*) how to efficiently
solve these tasks using LLMs and multi-agent collaboration.



Figure 1: Overview of AutoKaggle. AutoKaggle integrates a phase-based workflow with specialized agents (Reader, Planner, Developer, Reviewer, and Summarizer), iterative debugging and unit testing, a comprehensive machine learning tools library, and detailed reporting.

The core concept of AutoKaggle is phase-based multi-agent reasoning. This method leverages LLMs to reason and solve tasks within a structured workflow, addressing different facets of the data science process through the collaboration of multiple agents. AutoKaggle comprises two main components: a phase-based workflow and a multi-agent system, which complement each other, as shown in Figure 1.

134 **Phase-based Workflow.** The data science process is divided into six key phases: understanding the 135 background, preliminary exploratory data analysis, data cleaning, in-depth exploratory data anal-136 vsis, feature engineering, and model-building, -validation, and -prediction. Data cleaning, feature engineering, and model-building, -validation, and -prediction are fundamental processes required 137 for any data science competition. We designed two additional data analysis phases to provide essen-138 tial information and insights for data cleaning and feature engineering, respectively. Given that our 139 initial input is only an overview of a Kaggle data science competition and the raw dataset, we added 140 a background understanding phase to analyze various aspects of the competition background, objec-141 tives, file composition, and data overview from the raw input. This structured approach ensures that 142 all aspects of the problem are systematically and comprehensively addressed, with different phases 143 decoupled from each other. It allows thorough unit testing at each phase to ensure correctness and 144 prevent errors from propagating to subsequent phases. 145

Multi-agent System. The system consists of five specialized agents: Reader, Planner,
 Developer, Reviewer, and Summarizer. Each agent is designed to perform specific tasks
 within the workflow. They collaborate to analyze the problem, develop strategies, implement solutions, evaluate results, and generate comprehensive reports. Detailed setup and interaction processes of agents are described in Appendix D.1.



Figure 2: Iterative debugging and testing.

127

151 152

153

157

159

We summarize the pseudo-code of AutoKaggle in Algorithm 1. Let C represent the competition, Dthe dataset, and $\Phi = \{\phi_1, \phi_2, \dots, \phi_6\}$ the set of all phases in the competition workflow. For each phase ϕ_i , a specific set of agents \mathcal{A}_{ϕ_i} is assigned to perform various tasks. The key agents include Planner, Developer, Reviewer, and Summarizer.

166

- 167 168
- 169

2.2 DEVELOPMENT BASED ON ITERATIVE DEBUGGING AND TESTING

 In AutoKaggle, the Developer adopts a development approach based on iterative error correction and testing. It ensures the robustness and correctness of generated code through iterative execution, debugging, and testing.

Figure 2 shows the overall process of iterative debugging and testing. Specifically, the Developer first generates code based on the current state \mathbf{s}_t , the plan P_{ϕ_i} created by the Planner, and the historical context \mathcal{H} : $C_{\phi_i} = \text{GenerateCode}(\mathbf{s}_t, P_{\phi_i}, \mathcal{H})$. C_{ϕ_i} is the generated code for phase ϕ_i , and GenerateCode(\cdot) represents the code generation function executed by the Developer. The historical context \mathcal{H} includes previous phases' code, outputs, and other relevant information from other agents' activities.

After the initial code generation, it enters an iterative debugging and testing process. This process can be described by Algorithm 2.

182 Developer utilize three primary tools: code execution, code debugging, and unit testing.

(i) Code Execution. The Code Execution tool runs the generated code and captures any runtime errors. When an error is detected, the system restores a file to record the error messages.

(*ii*) Code Debugging. The Code Debugging tool analyzes error messages and attempts to fix the code. It utilizes error messages along with the current code and historical context to generate fixes: $C'_{\phi_i} = \text{DebugCode}(C_{\phi_i}, E_{\phi_i}, \mathcal{H}). C'_{\phi_i}$ is the debugged version of the code.

Following previous work (Tyen et al., 2024), we designed the debugging process into three main 189 steps: error localization, error correction, and merging of correct and corrected code segments. We 190 set a maximum of 5 attempts for the Developer to self-correct errors. Additionally, we've intro-191 duced an assistance mechanism. We record all error messages encountered during the debugging 192 process. When the number of correction attempts reaches 3, the Developer evaluates the feasibil-193 ity of continuing based on historical information. If past error messages are similar, it suggests that 194 the Developer might lack the ability to resolve this particular error, and continuing might lead to a loop. In such cases, we allow the Developer to exit the correction process and regenerate the 196 code from scratch. 197

(iii) Unit Testing. Unit testing runs predefined tests to ensure code meets requirements. For each phase ϕ_i , a set of unit tests T_{ϕ_i} is defined: $T_{\phi_i} = \{t_1, t_2, \dots, t_k\}$. The unit testing process can be represented as: $R_{\phi_i} = \text{ExecuteUnitTests}(C_{\phi_i}, T_{\phi_i})$. R_{ϕ_i} is the set of test results, with each result $r_j \in \{0, 1\}$ indicating whether the corresponding test passed (1) or failed (0).

In complex and accuracy-demanding tasks like Kaggle data science competitions, merely ensuring 202 that the code runs without errors is not enough. These competitions often involve intricate data 203 processing and sophisticated algorithms, where hidden logical errors can significantly affect the final 204 results. Therefore, it is necessary to design meticulous unit tests that not only verify the correctness 205 of the code but also ensure it meets the expected logical and performance standards. Otherwise, 206 hidden errors may accumulate through successive phases, making the completion of each subsequent 207 phase increasingly difficult. For example, unnoticed logical defects during the data cleaning phase 208 may lead to poor feature extraction, thereby affecting the model building in subsequent phases. 209

To mitigate these risks, unit tests for each phase must be carefully designed to cover a wide range of scenarios, including edge cases and potential failure points. This involves not only checking the correctness of the output but also ensuring that the intermediate steps conform to the expected logic. For instance, in the data cleaning phase, unit tests should verify whether missing values are handled correctly, outliers are appropriately managed, and data transformations are accurately applied.

By implementing comprehensive unit tests, we can catch and correct errors early in the development process, preventing them from propagating to later phases. This systematic testing approach ensures

that the code at each phase is not only error-free but also functionally correct and aligned with the overall project goals.

In conclusion, the iterative debugging and testing method employed by Developer ensures the generation of robust, error-free, and effective code for each phase of the competition. By employing advanced error handling, iterative debugging, and comprehensive unit testing, the system can adapt to various challenges and consistently produce high-quality code outputs.

222 223 224

2.3 MACHINE LEARNING TOOLS LIBRARY

225 Generating machine learning code from scratch using LLMs can be challenging due to the intri-226 cacies of various tasks. These models need to encompass specialized knowledge across a range 227 of processes, from data processing and feature engineering to model-building, -validation, and prediction. In many cases, leveraging expert-crafted machine learning tools is more efficient than 228 relying solely on LLM-generated code. This is because LLMs often lack domain-specific expertise, 229 potentially leading to suboptimal or inaccurate code. Furthermore, when tasked with complex oper-230 ations, the generated code may suffer from syntactical or logical errors, increasing the likelihood of 231 failures. 232

233 Our machine learning library is categorized into three core toolsets: data cleaning, feature engineer-234 ing, and model-building, -validation, and -prediction, each serving a specific role in the workflow. The data cleaning toolkit comprises seven tools, including FillMissingValues, RemoveColumns 235 WithMissingData, DetectAndHandleOutliersZscore, DetectAndHandleOutliersIqr, RemoveDupli-236 cates, ConvertDataTypes and FormatDatetime, all designed to ensure clean, consistent, and reli-237 able data preparation. The feature engineering module encompasses eleven tools aimed at enhanc-238 ing model performance, such as OneHotEncode, FrequencyEncode, CorrelationFeatureSelection, 239 and ScaleFeatures, employing various techniques like correlation analysis and feature scaling to 240 optimize data representation. The model-building, -validation, and -prediction category provides 241 TrainAndValidationAndSelectTheBestModel to support the full model development lifecycle, in-242 cluding model selection, training, evaluation, prediction, ensemble integration, and hyperparame-243 ter optimization, facilitating robust model deployment and effective performance. Each tool comes 244 with comprehensive explanations, input/output specifications, anomaly detection, and error handling 245 guidance.

246 This comprehensive library is crucial for efficient multi-agent collaboration in tackling complex 247 Kaggle competitions. Each tool provides standardized, reliable functionality, enabling AutoKaggle 248 to seamlessly share and process data, enhance feature quality, and optimize model performance, 249 ultimately improving overall workflow efficiency and ensuring coordinated, high-quality solutions 250 in a competitive environment. Moreover, our machine learning library reduces the burden on AutoKaggle in detailed programming tasks, enabling them to focus more on higher-level task planning 251 and code design. This shift of focus allows AutoKaggle to navigate complex tasks more effectively, 252 ultimately improving their overall performance. More details of our machine learning tools can be 253 found in Appendix D.3. 254

255 256

257

3 EXPERIMENTS

258 3.1 EXPERIMENTAL SETUP

259 Task Selection. We select eight Kaggle competitions that predominantly use tabular datasets, focus-260 ing on classification and regression tasks. These competitions are categorized into two types: *classic* 261 Kaggle and Recent Kaggle. Classic Kaggle competitions are those that begin before October 2023 262 with at least 500 participants, whereas Recent Kaggle competitions begin in 2024 or later. As our 263 analysis relies on GPT-40, which is trained on data available until October 2023, it includes most 264 of the Classic Kaggle competitions. To evaluate the generalization capabilities of AutoKaggle, we 265 therefore focus on competitions initiated after 2024. Additionally, we classify these competitions 266 into three difficulty levels: easy, medium, and hard. For each dataset, we access the corresponding 267 competition's homepage on Kaggle, extract content from the overview and data description sections, and compile this information into a file named overview.txt. This file, along with the original com-268 petition data files, forms the primary input for AutoKaggle. More details of our datasets can be 269 found in Appendix C.

270 Notably, we do not incorporate the nine tabular datasets from MLE-Bench (Hong et al., 2024) due 271 to their substantial size, which would significantly increase computational runtime. Resource lim-272 itations prevent us from adhering to MLE-Bench's experimental setup, which specifies a 24-hour 273 participation window per agent and a 9-hour code execution timeout.

274 275

276

277

Table 1: Made submission, valid submission and comprehensive score on 8 Kaggle tasks. Each experiment is repeated with 5 trials. The best performances on individual tasks are underlined, and the best performances across all tasks are bolded.

Metric	Setting / Task	Classic			Recent					
	Setting / Tusk	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	Task 7	Task 8	Avg.
Made Submission	AutoKaggle gpt-4o AutoKaggle o1-mini AIDE gpt-4o	$\frac{\underline{1}}{\underline{1}}$	$\frac{0.80}{0.60}$ 0.40	$\frac{0.80}{0.60}$ 0.20	$\begin{array}{c} \frac{1}{1} \\ 0.60 \end{array}$	$0.80 \\ 0.60 \\ \underline{1}$	$\frac{0.80}{0.80}\\ \frac{0.80}{0.80}$	<u>0.80</u> 0.60 <u>0.80</u>	$\begin{array}{c} \underline{0.80}\\ 0.60\\ 0 \end{array}$	0.85 0.73 0.60
Valid Submission	AutoKaggle gpt-4o AutoKaggle o1-mini AIDE gpt-4o	$\frac{\underline{1}}{\underline{1}}$	<u>0.80</u> 0.60 0.40	<u>0.80</u> 0.60 0.20	$\frac{1}{1}$ 0.40	$\begin{array}{c} \underline{0.80}\\ 0.60\\ \underline{1} \end{array}$	0.60 0.60 <u>0.80</u>	<u>0.80</u> 0.60 <u>0.80</u>		0.83 0.70 0.58
Comprehensive Score	AutoKaggle gpt-40 AutoKaggle o1-mini AIDE gpt-40	<u>0.888</u> 0.879 0.872	<u>0.786</u> 0.680 0.597	<u>0.831</u> 0.729 0.542	0.862 <u>0.863</u> 0.561	0.810 0.709 <u>0.918</u>	0.728 0.735 <u>0.793</u>	0.848 0.742 0.848	$\frac{0.812}{0.735}$ 0	0.821 0.759 0.641

287 Evaluation metric. We evaluate the capability of the AutoKaggle from four perspectives: Made 288 Submission, Valid Submission, Average Normalized Performance Score and Comprehensive Score. 289 The first two metrics refer to MLE-bench and are primarily used to assess the ability to generate a 290 submission.csv file. The last two metrics come from Data Interpreter (Chan et al., 2024), we made 291 modifications to adapt them to the evaluation of our framework.

292 (i) Made Submission (MS). Made Submission refers to the percentage of times a submission.csv 293 file is generated. 294

(ii) Valid Submission (VS). Valid Submission indicates the percentage of those submission.csv files 295 that are valid—meaning they can be successfully submitted to the Kaggle website, produce results 296 without errors, and have no issues related to data scale or category mismatches. 297

298 (iii) Comprehensive Score (CS). In the evaluations, performance metrics are divided into two cate-299 gories: bounded metrics, which range from 0 to 1 where higher values indicate better performance, 300 and unbounded metrics, where lower values denote superior performance. To normalize these different types of metrics, we utilize the normalized performance score (NPS), defined as follows: 301

$$NPS = \begin{cases} \frac{1}{1+s}, & \text{if } s \text{ is smaller the better} \\ s, & \text{otherwise.} \end{cases}$$
(1)

For multiple trials of a task, we calculate the Average Normalized Performance Score (ANPS) as the average of the successful attempts:

$$ANPS = \frac{1}{T_s} \sum_{t=1}^{T_s} NPS_t$$
(2)

309 310 311

312 313

314

307

308

> Table 2: Ablation study on machine learning tools. Evaluated with completion rate and comprehensive score. Best performance are underlined.

		Task 1	Task 2	Task 3	Task 5	Avg.
	No Tools	0.80	0.60	0.50	0.40	0.58
VC	DC Tools	0.80	0.70	1.00	1.00	0.88
V 3	DC & FE Tools	0.80	0.60	$\overline{0.60}$	0.60	0.65
	All Tools	<u>1.00</u>	<u>0.80</u>	0.80	0.80	0.85
	No Tools	0.781	0.697	0.666	0.602	0.68
~~	DC Tools	0.781	0.721	0.928	0.909	0.83
CS	DC & FE Tools	0.787	0.684	0.735	0.713	0.73
	All Tools	0.888	0.786	0.831	0.810	0.82

324 where T_s represents the total number of successful attempts for a task, and NPS_t is the NPS value 325 for the *t*-th attempt. 326

To comprehensively evaluate both the pass rate and the average performance, we define the Com-327 prehensive Score (CS) as the average of VS and ANPS: 328

$$CS = 0.5 \times VS + 0.5 \times ANPS$$
(3)

Experiment Details. We evaluated AutoKaggle's performance based on both GPT-40 and 01-mini models. Notably, different models were assigned to specific agents based on their functional re-332 quirements. The Reader, Reviewer, and Summarizer, which perform tasks requiring min-333 imal logical reasoning and coding capabilities, were implemented using the GPT-4o-mini model. 334 The Planner, responsible for task decomposition and planning that demands sophisticated logi-335 cal reasoning, operates on either the GPT-40 or o1-mini model. Although the Developer's tasks 336 traditionally necessitate advanced logical reasoning and coding skills, the Planner's effective task 337 decomposition methodology has moderated these requirements, therefore it is based on GPT-40 338 model.

339 In our experiments, Each task undergoes five trials, with each phase in the workflow allowing for a 340 maximum of three iterations. During an iteration, the Developer may debug the code up to five 341 times. If unsuccessful, they proceed with the same phase, deriving insights and adjusting strategies 342 based on previous attempts. Failure to resolve issues after three iterations is considered a definitive 343 failure. 344

Baseline. We employ AIDE (Schmidt et al., 2024) as our baseline, which is the best-performing 345 framework in MLE-bench evaluation results. We use AIDE's default settings, only modifying 346 agent.base.model to the GPT-40 model. 347

MAIN RESULTS 3.2

329 330

331

348

349 350

351

352

The comprehensive performance of AutoKaggle across 8 Kaggle data science competitions is presented in Table 1. In order to facilitate understanding, we uniformly name the eight tasks as task 1-8. The real task names and detailed dataset information are available in Appendix C.



Made submission and Valid submission. We first evaluated the success rate of valid submis-376 sion.csv file generation across different experimental configurations. The AutoKaggle framework, 377 implemented with GPT-40, demonstrated superior performance with an average valid submission



389

Figure 4: Left. Debugging time and Right. Average performance in competitions.

rate of 83% across all 8 Kaggle tasks, surpassing the AIDE framework by 28%. These results underscore the robustness of our framework in executing comprehensive data science workflows. While
the AIDE framework successfully processed Tasks 1-7, which involved single-variable classification or regression on tabular data, it failed to generate valid submissions for Task 8, a multi-variable
classification problem. This differential performance demonstrates our framework's versatility in
handling diverse tabular data tasks.

396 Another interesting observation is that within the AutoKaggle framework, the GPT-40 model 397 achieved better results than the o1-mini model, despite the latter's purported superior reasoning 398 capabilities. This performance difference emerged solely from varying the model used in the Planner component. We hypothesize that this counterintuitive result stems from o1-mini's tendency 399 toward excessive planning complexity, which proves disadvantageous in our streamlined, phase-400 based workflow architecture. This same consideration influenced our decision to maintain GPT-40 401 as the Developer's base model, as our experiments indicated that an ol-mini-based Developer 402 would significantly increase code verbosity, expanding 100-line solutions to approximately 500 lines 403 through the introduction of superfluous components such as logging systems. 404

Comprehensive Score. Subsequently, we compared the overall performance of different settings across 8 Kaggle tasks. AutoKaggle with GPT-40 achieved the highest comprehensive score in 5 tasks and demonstrated the best overall performance. Figure 3 illustrates the comparison of different settings based on the average normalized performance score metric, where AutoKaggle with ol-mini achieved the highest overall score. This indicates that although the ol-mini-based Planner generated overly complex plans that increased development difficulty, successfully executing these plans according to specifications led to superior performance outcomes.

411 412

413 3.3 ABLATION STUDY

Apart from the modules involved in the ablation study, all other experimental settings are identical to those in the formal experiment.

416 Study on Machine Learning Tools. To evaluate the effectiveness of the machine learning tools 417 module and the impact of tools across different phases on the results, we conduct ablation exper-418 iments. We begin without any tools and progressively add them at each phase until all machine 419 learning tools are implemented. The results are presented in Table 2. Notably, the completion rate 420 increases by 30% with the use of data cleaning phase tools, and by 27.5% when all tools are utilized, 421 compared to the scenario with no tools. However, the completion rate exhibits a decline during the 422 feature engineering phase, particularly in the house prices and academic success competitions. This 423 decline can be attributed to the relatively large number of features involved, alongside the complexity and high encapsulation of the tools used in this phase, which necessitate the addition and 424 removal of features, thereby complicating their usage. Furthermore, this complexity poses chal-425 lenges for Developers in debugging erroneous code. As illustrated in Figure 4 (a), the frequency 426 of debugging instances is greater when employing tools from the feature engineering phase. 427

Figure 4 (b) provides a clearer comparison, demonstrating that while the best normalized performance scores across four scenarios are similar, the completion rate significantly increases with the use of the tool. This suggests that although the machine learning tool library we develop does not substantially elevate the solution's upper limit, it functions as a more stable tool that enhances AutoKaggle's completion rate. This outcome aligns with expectations, as the machine learning tool

library is a redevelopment based on widely used libraries such as pandas and scikit-learn. It does not introduce new functionalities but instead combines and re-packages existing ones, incorporating error handling and manual testing to ensure compatibility with our framework.

Study on Unit Tests. To evaluate the effectiveness of the unit tests module, we conduct ablation experiments. The results are presented in Table 3. In the absence of unit tests, the completion rate significantly decreases, making it nearly impossible to complete the tasks. This emphasizes that for tasks like data science, which demand high levels of precision and logic, it is not enough for each phase of the code to merely execute without errors. Comprehensive unit testing is required to ensure that the code is logical and achieves the objectives of each phase.

Study on Debugging Times. We conduct ablation experiments to investigate the impact of the num-ber of allowed debugging times on the results. The experimental setup permits five code debugging attempts within each phase, with each phase being executable up to three times. Consequently, we analyze scenarios with allowable corrections set at 0, 5, and 10. The results are shown in Figure 5. It can be observed that when AutoKaggle is required to pass without any errors, there is only one successful record on the Titanic task. Allowing five debugging attempts significantly improves the completion rate, and further increases in allowable debugging attempts lead to rises in all metrics. This demonstrates the efficacy of our code debugging module. However, the performance when the number of allowable debugging attempts is set to 10 and 15, suggesting that the agent's self-correction abilities are limited. There are complex errors that it cannot resolve independently, and further increasing the number of allowable debugging attempts does not address these errors.See more details in Section B.



Figure 5: Comprehensive Score across different debugging times.

Study on Competition Date. To further evaluate the generalization capabilities of our AutoKaggle framework, we conducted an analysis stratified by competition dates. Tasks 1-4 corresponded to competitions potentially included in the training data of models such as GPT-40 and O1-mini, while tasks 5-8 were derived from competitions launched in the current year. This temporal stratification enabled us to assess the framework's performance on out-of-distribution tasks. For classic Kaggle tasks, AutoKaggle with GPT-40 achieved a valid submission rate of 0.90 and a comprehensive score of 0.842. On recent tasks, these metrics were 0.75 and 0.800 respectively, demonstrating only marginal performance degradation. These results indicate that our task decoupling approach and

Table 3: Ablation study on unit tests. Better performance are underlined.

r						
		Task 1	Task 2	Task 3	Task 5	Avg.
CR	w/o Unit Tests	0.20	0	0.20	0	0.10
		<u>1.00</u>	<u> </u>	<u> </u>	<u>0.80</u> 	<u>0.83</u>
CS	w/o Unit Tests	0.478	0	0.482	0	0.240
CD	w/ Unit Tests	0.888	0.831	0.786	0.810	0.829

486 predefined execution pathways enable effective handling of novel competitions, even in scenarios 487 where the underlying model lacks prior exposure to the domain. 488

RELATED WORK 4

490

489

491 A concise framework of agents consists of brain, perception, and action modules (Xi et al., 2023). 492 The perception module processes external information, the brain plans based on that information, 493 and the action module executes these plans (Xi et al., 2023; Zhou et al., 2023). LLMs, acting as brain 494 modules, exhibit impressive zero-shot abilities and are applied in fields like data science and music 495 composition (Brown et al., 2020; Hong et al., 2024; Deng et al., 2024). While the chain-of-thought 496 method enhances reasoning (Wei et al., 2023), it still faces challenges related to hallucinations and 497 unfaithfulness (Turpin et al., 2023), potentially due to internal representations (Yao et al., 2023). 498 The ReAct paradigm addresses this by integrating thoughts and actions, refining outputs through 499 interaction with external environments (Yao et al., 2023; Madaan et al., 2023; Shinn et al., 2023;

500 Zhou et al., 2024).

501 While an individual agent can achieve basic natural language processing (NLP) tasks, real-world 502 tasks have higher complexities. In human societies, people chunk complex tasks into simple sub-503 tasks that different people can easily handle. Inspired by this division of labor principle, multi-agent 504 systems enhance performance (Talebirad & Nadiri, 2023) using cooperative interactions (Xi et al., 505 2023; Li et al., 2023) to achieve shared goals. Another interaction method is adversarial interactions 506 (Lewis et al., 2017), where several agents compete with each other for better results, or one agent 507 critiques and reviews the generation of another agent (Gou et al., 2024).

508 In order to address the well-defined requirements of data science tasks, a feasible approach is to de-509 sign hierarchical systems (Hong et al., 2024; Zhang et al., 2024b; Chi et al., 2024) to complete tasks 510 such as task understanding, feature engineering, and model building. In each hierarchy, separately 511 design two agents for the code planning and code generation respectively (Hong et al., 2024), and use 512 unit tests (Zhang et al., 2024b) to verify the quality of code generation. Beyond self-debugging by 513 autonomous multi-agents, human-in-the-loop (Hong et al., 2024; Zhang et al., 2024b) mechanisms 514 also provide oversight and corrections to LLM outputs, reducing hallucinations in each hierarchy. 515 Tang et al. (2024) introduces ML-Bench, a benchmark for language agents for machine learning tasks. 516

517 In summary, multi-agent systems and LLM-based agents have demonstrated significant potential 518 across domains such as NLP and data science. While individual agents excel in basic tasks, inte-519 grating multiple agents is crucial for tackling complex real-world challenges. By combining task-520 specific agents with human-in-the-loop mechanisms and unit testing, these systems improve code quality and address issues like hallucinations. Our framework, AutoKaggle, advances these efforts 521 by integrating LLM-based reasoning with multi-agent collaboration, ensuring adaptability, correct-522 ness, and user control in data science competitions. 523

524

5 CONCLUSION

525 526 527

In this paper, we introduce AutoKaggle, a robust framework designed to leverage phase-based work-528 flows and multi-agent collaboration for solving complex Kaggle data science competitions. Au-529 toKaggle employs an iterative development process, incorporating thorough code debugging, unit 530 testing, and a specialized machine learning tools library to address the intricate requirements of data 531 science tasks. Our framework enhances reliability and automation in managing sophisticated data workflows, while maintaining user control through customizable processes. Extensive evaluations 532 across various Kaggle competitions demonstrate AutoKaggle's effectiveness, marking a significant 533 advancement in AI-assisted data science problem-solving and expanding the capabilities of LLM-534 based systems in tackling real-world challenges. 535

536

REFERENCES 538

Ryan Holbrook Addison Howard, Ashley Chow. Spaceship titanic, 2022. URL https:// kaggle.com/competitions/spaceship-titanic.

540 DataCanary Anna Montoya. House prices advanced regression tech-541 niques, 2016. URL https://kaggle.com/competitions/ 542 house-prices-advanced-regression-techniques. 543 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-544 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. 546 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz 547 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec 548 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL 549 https://arxiv.org/abs/2005.14165. 550 551 Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio 552 Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, et al. Mle-bench: Evaluating machine learning agents on machine learning engineering. arXiv preprint arXiv:2410.07095, 2024. 553 554 Yizhou Chi, Yizhang Lin, Sirui Hong, Duyi Pan, Yaying Fei, Guanghao Mei, Bangbang Liu, Tianqi 555 Pang, Jacky Kwok, Ceyao Zhang, Bang Liu, and Chenglin Wu. Sela: Tree-search enhanced 556 llm agents for automated machine learning, 2024. URL https://arxiv.org/abs/2410. 17238. 558 559 Will Cukierski. Titanic - machine learning from disaster, 2012. URL https://kaggle.com/ competitions/titanic. 561 Qixin Deng, Qikai Yang, Ruibin Yuan, Yipeng Huang, Yi Wang, Xubo Liu, Zeyue Tian, Jiahao 562 Pan, Ge Zhang, Hanfeng Lin, Yizhi Li, Yinghao Ma, Jie Fu, Chenghua Lin, Emmanouil Benetos, 563 Wenwu Wang, Guangyu Xia, Wei Xue, and Yike Guo. Composerx: Multi-agent symbolic music 564 composition with llms, 2024. URL https://arxiv.org/abs/2404.18081. 565 566 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 567 Critic: Large language models can self-correct with tool-interactive critiquing, 2024. URL 568 https://arxiv.org/abs/2305.11738. 569 Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. Ds-agent: Auto-570 mated data science by empowering large language models with case-based reasoning, 2024. URL 571 https://arxiv.org/abs/2402.17453. 572 573 Md Mahadi Hassan, Alex Knipper, and Shubhra Kanti Karmaker Santu. Chatgpt as your personal 574 data scientist, 2023. URL https://arxiv.org/abs/2305.13657. 575 Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Danyang Li, Jiaqi Chen, Jiayi 576 Zhang, Jinlin Wang, Li Zhang, Lingyao Zhang, Min Yang, Mingchen Zhuge, Taicheng Guo, Tuo 577 Zhou, Wei Tao, Wenyi Wang, Xiangru Tang, Xiangtao Lu, Xiawu Zheng, Xinbing Liang, Yaying 578 Fei, Yuheng Cheng, Zongze Xu, and Chenglin Wu. Data interpreter: An llm agent for data 579 science, 2024. URL https://arxiv.org/abs/2402.18679. 580 581 Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su, 582 Jingjing Xu, Ming Zhu, Yao Cheng, Jianbo Yuan, Jiwei Li, Kun Kuang, Yang Yang, Hongxia 583 Yang, and Fei Wu. Infiagent-dabench: Evaluating agents on data analysis tasks, 2024. URL 584 https://arxiv.org/abs/2401.05507. 585 Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming 586 Zhang, Xinya Du, and Dong Yu. Dsbench: How far are data science agents to becoming data science experts?, 2024. URL https://arxiv.org/abs/2409.07703. 588 589 Wendy Kan. Ghouls, goblins, and ghosts... boo!, 2016. URL https://kaggle.com/ competitions/ghouls-goblins-and-ghosts-boo. 591 Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? end-592 to-end learning for negotiation dialogues, 2017. URL https://arxiv.org/abs/1706. 593 05125.

- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem.
 Camel: Communicative agents for "mind" exploration of large language model society, 2023.
 URL https://arxiv.org/abs/2303.17760.
- Killian Lucas. GitHub KillianLucas/open-interpreter: A natural language interface for computers
 github.com. https://github.com/KillianLucas/open-interpreter, 2023.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad
 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self refine: Iterative refinement with self-feedback, 2023. URL https://arxiv.org/abs/
 2303.17651.
- 605 OpenAI. Chatgpt: Optimizing language models for dialogue. https://openai.com/blog/ chatgpt, 2022.
- OpenAI. Gpt-4 technical report. https://arxiv.org/abs/2303.08774, 2023.
- Dominik Schmidt, Zhengyao Jiang, and Yuxiang Wu. Introducing Weco AIDE. https://www.
 weco.ai/blog/technical-report, April 2024.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL https://arxiv.org/abs/2303.11366.
- Yashar Talebirad and Amirhossein Nadiri. Multi-agent collaboration: Harnessing the power of
 intelligent llm agents, 2023. URL https://arxiv.org/abs/2306.03314.
- Kiangru Tang, Yuliang Liu, Zefan Cai, Yanjun Shao, Junjie Lu, Yichi Zhang, Zexuan Deng, Helan Hu, Kaikai An, Ruijun Huang, Shuzheng Si, Sheng Chen, Haozhe Zhao, Liang Chen, Yan Wang, Tianyu Liu, Zhiwei Jiang, Baobao Chang, Yin Fang, Yujia Qin, Wangchunshu Zhou, Yilun Zhao, Arman Cohan, and Mark Gerstein. Ml-bench: Evaluating large language models and agents for machine learning tasks on repository-level code, 2024. URL https://arxiv.org/abs/2311.09835.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting, 2023. URL https://arxiv.org/abs/2305.04388.
- Gladys Tyen, Hassan Mansoor, Victor Cărbune, Peter Chen, and Tony Mak. Llms cannot find
 reasoning errors, but can correct them given the error location, 2024. URL https://arxiv.
 org/abs/2311.08516.
- Ashley Chow Walter Reade. Binary classification with a bank churn dataset, 2024a. URL https:
 //kaggle.com/competitions/playground-series-s4e1.
- Ashley Chow Walter Reade. Multi-class prediction of obesity risk, 2024b. URL https: //kaggle.com/competitions/playground-series-s4e2.
- Ashley Chow Walter Reade. Steel plate defect prediction, 2024c. URL https://kaggle.com/
 competitions/playground-series-s4e3.
- Ashley Chow Walter Reade. Classification with an academic success dataset, 2024d. URL https: //kaggle.com/competitions/playground-series-s4e6.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc
 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models,
 2023. URL https://arxiv.org/abs/2201.11903.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao
 Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou,
 Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey,
 2023. URL https://arxiv.org/abs/2309.07864.

648 649 650	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023. URL https://arxiv. org/abs/2210.03629.
652 653 654	Lei Zhang, Yuge Zhang, Kan Ren, Dongsheng Li, and Yuqing Yang. Mlcopilot: Unleashing the power of large language models in solving machine learning tasks, 2024a. URL https://arxiv.org/abs/2304.14979.
655 656	Yaolun Zhang, Yinxu Pan, Yudong Wang, and Jie Cai. Pybench: Evaluating llm agent on various real-world coding tasks, 2024b. URL https://arxiv.org/abs/2407.16732.
658 659	Yuge Zhang, Qiyang Jiang, Xingyu Han, Nan Chen, Yuqing Yang, and Kan Ren. Benchmarking data science agents, 2024c. URL https://arxiv.org/abs/2402.17168.
660 661 662 663 664	Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Xiangru Tang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. Agents: An open-source framework for autonomous language agents, 2023. URL https://arxiv.org/abs/2309.07870.
665 666 667 668	Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, Huajun Chen, and Yuchen Eleanor Jiang. Sym- bolic learning enables self-evolving agents, 2024. URL https://arxiv.org/abs/2406. 18532.
669	
671	
670	
672	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
700	
701	

702 A ALGORITHM

```
704
705
            Algorithm 1: AutoKaggle Workflow
706
            Input : Competition C, Dataset D
707
            Output: Solution S, Comprehensive report \mathcal{R}
         1 Initialize state s_0 with first phase \phi_1: "Understand Background";
708
         2 t \leftarrow 0;
709
         \bullet \Phi \leftarrow \{\phi_1, \phi_2, ..., \phi_6\};
                                                                                                    /* Set of all phases */
710
         4 Define \mathcal{A}_{\phi} for each \phi \in \Phi;
                                                                                           /* Agents for each phase */
711
         5 do
712
                  \mathbf{s}_t \leftarrow \text{GetCurrentState}();
         6
713
                  \phi_{\text{current}} \leftarrow \text{GetCurrentPhase}(\Phi);
         7
714
                  \mathcal{A}_t \leftarrow \mathcal{A}_{\phi_{\text{current}}};
         8
715
                  for a \in \mathcal{A}_t do
         9
716
                       if a is Planner then
         10
717
                            r_a \leftarrow a.execute(\mathbf{s}_t);
         11
                             \mathbf{s}_t \leftarrow \text{UpdateState}(\mathbf{s}_t, r_a);
718
         12
719
                            if UserInteractionEnabled() then
         13
                              \mathbf{s}_t \leftarrow \text{UserReview}(\mathbf{s}_t);
                                                                                                      /* User Review plan */
720
         14
721
                       else if a is Developer then
         15
722
                            r_a \leftarrow a.execute(\mathbf{s}_t);
         16
723
                             \mathbf{s}_t \leftarrow \text{UpdateState}(\mathbf{s}_t, r_a);
         17
724
                             if NoErrors(r_a) then
         18
725
                                  T \leftarrow \text{ExecuteUnitTests}(\phi_{\text{current}});
         19
726
                                  if \neg PassTests(T) then
         20
                                       \mathbf{s}_t \leftarrow \text{Debug}(\mathbf{s}_t);
727
         21
728
                       else
        22
729
                             r_a \leftarrow a.\text{execute}(\mathbf{s}_t);
        23
730
                             \mathbf{s}_t \leftarrow \text{UpdateState}(\mathbf{s}_t, r_a);
         24
731
732
                  if AllAgentsCompleted(A_t) and PassTests(T) then
        25
733
                   \phi_{\text{current}} \leftarrow \text{NextPhase}(\Phi);
        26
734
        27
                 t \leftarrow t + 1;
735
        28 while \exists \phi \in \Phi : not completed(\phi);
736
        29 S \leftarrow \text{ExtractSolution}(\mathbf{s}_t);
737
        30 \mathcal{R} \leftarrow \text{GenerateReport}(\mathbf{s}_t);
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
```

756 Algorithm 2: Development based on Iterative Debugging and Testing **Input** : Initial code C_{ϕ_i} , Current state \mathbf{s}_t , Plan P_{ϕ_i} , Historical context \mathcal{H} , Maximum tries 758 max_tries , Error threshold threshold759 **Output:** Debugged and tested code C'_{ϕ_i} , Execution flag *execution_flag* 760 1 round $\leftarrow 0$; 761 2 $error_flag \leftarrow false;$ 762 s execution_flag \leftarrow true; 4 retry_flag \leftarrow false; 764 5 error_history $\leftarrow \emptyset$; 765 while round < max_tries do 6 766 if round = 0 or $retry_flag$ then 7 767 $C_{\phi_i} \leftarrow \text{GenerateCode}(\mathbf{s}_t, P_{\phi_i}, \mathcal{H});$ 8 768 *error_history* $\leftarrow \emptyset$; 9 $retry_flag \leftarrow false;$ 769 10 $error_flag, E_{\phi_i} \leftarrow \operatorname{RunCode}(C_{\phi_i});$ 770 11 if error_flag then 12 771 $error_history \leftarrow error_history \cup \{E_{\phi_i}\};$ 13 if $|error_history| >= threshold$ then 14 15 $retry_flag \leftarrow EvaluateRetry(error_history);$ 774 if *retry_flag* then 16 775 continue; 17 776 $C_{\phi_i} \leftarrow \text{DebugCode}(C_{\phi_i}, E_{\phi_i}, \mathcal{H});$ 777 18 778 19 else $R_{\phi_i} \leftarrow \text{ExecuteUnitTests}(C_{\phi_i}, T_{\phi_i});$ 20 780 21 if $\exists r_j \in R_{\phi_i} : r_j = 0$ then $\check{C}_{\phi_i} \leftarrow DebugTestFailures(C_{\phi_i}, R_{\phi_i}, \mathcal{H});$ 781 22 782 else 23 783 execution_flag \leftarrow true; 24 784 break; 25 785 $round \leftarrow round + 1;$ 26 786 if $round = max_tries$ then 27 787 execution_flag \leftarrow false; 28 788 789 29 return C_{ϕ_i} , execution_flag 790

791 792

793

794

B ERROR ANALYSIS

In each subtask phase of AutoKaggle, errors may occur, with data cleaning and feature engineering experiencing the highest error rates at 25% and 22.5%, respectively. Notably, failures during the feature engineering phase result in direct competition failures in 31.25% of cases.

798 In the context of the proposed AutoKaggle framework, which aims to assist data scientists in solv-799 ing complex tabular data challenges through a collaborative multi-agent system, Table 4 provides 800 an overview of the different types of errors encountered during the iterative development process. 801 AutoKaggle's workflow includes code execution, debugging, and comprehensive unit testing, and 802 the listed errors are indicative of the various challenges encountered while automating these stages. The most frequently observed errors are Value Errors (49 occurrences), related to mismatched input 804 types or ranges, and Key Errors (44 occurrences), resulting from attempts to access non-existent dic-805 tionary keys. Additionally, Type Errors (25 occurrences) and Model Errors (8 occurrences) highlight operational issues due to data type mismatches or incorrect model configurations, respectively. The table also details other errors such as Timeout, FileNotFound, and Index Errors, each contributing to the debugging process. Understanding these error types is crucial for improving AutoKaggle's 808 robustness and aligning automated workflows with human interventions, ultimately enhancing pro-809 ductivity in typical data science pipelines.

In addition, we provide a detailed debugging process for developers. Below, we illustrate this using a FileNotFoundError as an example of the debugging workflow:

- Error Localization: The developer initially encounters issues executing a Python script involving file-saving operations with libraries like Matplotlib and Pandas. The specific error, FileNotFoundError, is traced to nonexistent directories or incorrect file paths. Through an iterative analysis, the problematic sections of the code are identified, focusing on the need to properly manage directory paths and handle filenames.
- Error Correction: To address these issues, several modifications are suggested. First, the importance of ensuring that directories exist before performing file operations is highlighted by incorporating os.makedirs to create any missing directories. Additionally, a filename sanitization approach is recommended to prevent errors related to invalid characters in file paths. A custom sanitize_filename function is introduced to ensure filenames contain only valid characters, thereby avoiding issues caused by special symbols or whitespace.
 - Merging Correct and Corrected Code Segments: The final step involves merging the corrected segments back into the original code to create a seamless and robust solution. The revised script includes improvements such as verifying directory existence, creating necessary directories, and applying filename sanitization to ensure compatibility across different operating systems. The corrected code is delivered with a focus on enhancing reliability, particularly in file-saving processes, making it resilient against common pitfalls like missing directories or invalid filenames.
- C DETAILED DATASET DESCRIPTION

Here is the detailed description of our dataset. Note that we use task labels to represent the different datasets. Task 1 refers to Titanic (Cukierski, 2012), Task 2 refers to Spaceship Titanic (Addison Howard, 2022), Task 3 refers to House Prices (Anna Montoya, 2016), Task 4 refers to Monsters (Kan, 2016), Task 5 refers to Academic Success (Walter Reade, 2024d), Task 6 refers to Bank Churn (Walter Reade, 2024a), Task 7 refers to Obesity Risk (Walter Reade, 2024b), and Task 8 refers to Plate Defect (Walter Reade, 2024c).

Our framework deliberately avoids selecting competitions with excessively large datasets. The reason for this is that larger datasets significantly extend the experimental runtime, making it impractical to dedicate a machine to a single experiment for such prolonged periods.

Table 4: Error Types of AutoKaggle in the Problem-Solving Stage

Error Type (Count)	Description
Value Error (49)	Fail to match the expected type or range of the input values
Key Error (44)	Attempt to access a dictionary element using a key that does not exist
File Error (8)	Attempt to access a file that does not exist in the specified location
Model Error (8)	Incorrect setup in the parameters or structure of a model, leading to opera- tional failures
Type Error (25)	Mismatch between expected and actual data type, leading to operational failure
Timeout Error (6)	Failure to complete a process within the allocated time period
Index Error (3)	Attempt to access an element at an index that is outside the range of a list or array
Assertion Error (1)	An assertion condition in the code is not met, indicating an unmet expected constraint
Name Error (2)	Use of an undeclared variable that is not recognized by the system
Attribute Error (2)	Attempt to access an attribute or method that does not exist for an object
Indentation Error (1)	Incorrect indentation disrupts code structure, preventing proper parsing

First, we intentionally avoided selecting competitions with datasets that were too large, as larger datasets can significantly extend the experimental runtime, making it impractical to use a single machine for extended experiments. Second, we adhered to real-world competition settings by gen-erating submission files and submitting them manually for evaluation. Simply splitting the training data would result in a test set with a distribution very similar to the training data, which could in-flate performance metrics-similar to the difference often seen between validation scores and real test scores. Third, our dataset clearly identifies the contest type, i.e., tabular data. Fourth, since datasets for large language modeling include publicly available Kaggle contests, we selected only those released after 2024. Our framework requires agents to independently interpret contest tasks, understand the data, and determine appropriate optimization strategies without relying on predefined guidance."

Table 5: Selected Kaggle tasks. For each task, we show its number, category, difficulty level, numberof teams, train size and test size in dataset.

Category	No.	Task Name	Task	Level	Teams	Train	Test
	1	Titanic	Classification	Medium	13994	891	418
Classia	2	Spaceship Titanic	Classification	Easy	1720	8693	4277
Classic	3	House Prices	Regression	Medium	4383	1460	1459
	4	Monsters	Classification	Easy	763	371	529
	5	Academic Success	Regression	Medium	2684	76.5K	51K
Decent	6	Bank Churn	Regression	Easy	3632	165K	110K
Ketellt	7	Obesity Risk	Classification	Easy	3587	20.8K	13.8K
	8	Plate Defect	Regression	Medium	2199	19.2K	12.8K

D IMPLEMENTATION DETAILS

D.1 AGENT DETAILS

D.1.1 AGENT BASE

The base agent is a father class of other agents (Reader, Planner, Developer, Reviewer, and Summarizer) in the AutoKaggle. This agent can act with various tools for tasks related to data analysis, model evaluation, and document retrieval etc.

D.1.2 READER

Reader is designed for reading documents and summarizing information. It processes overview.txt in each competition, subsequently providing a well-organized summary of the competition's back-ground

Prompt of Agent Reader / Task Prompt							
Role: reading documents and summarizing information							
Description: The Reader only appears in the Understand							
Background phase, it reads the overview.txt file of the							
Kaggle competition, the sample data of both training and							
testing sets and summarizes it into a clearly structured							
competition_info.txt in markdown format.							



D.1.3 PLANNER

955

956 957

958

959

964 965

966

967

968

969

970

971

Planner is designed for creating task plans and roadmaps. The agent's main function is to structure and organize tasks into executable plans, primarily by leveraging available tools and previously generated reports.

Prompt of Agent Planner / Task Prompt

Role: creating task plans and roadmaps Description: In the first execution, the Planner collects the competition information, the current state, and the user's rules to generate a new plan. This generation involves several rounds of interaction with a LLM to gather task details, reorganize data into structured formats (Markdown and JSON), and finalize a plan.

972 Prompt of Agent Planner / Task Prompt 973 974 # CONTEXT # 975 {phases_in_context} 976 Currently, I am at phase: {phase_name}. 977 ############# 978 # INFORMATION # {background_info} 979 980 {state_info} 981 ############# 982 # NOTE # ## PLANNING GUIDELINES ## 983 1. Limit the plan to a MAXIMUM of FOUR tasks. 984 2. Provide clear methods and constraints for each task. 3. Focus on critical steps specific to the current phase. 985 4. Prioritize methods and values mentioned in USER RULES. 986 5. Offer detailed plans without writing actual code. 6. ONLY focus on tasks relevant to this phase, avoiding those belonging to other 987 phases. For example, during the in-depth EDA phase 988 - you CAN perform detailed univariate analysis on KEY features. - you CAN NOT modify any feature or modify data. 989 990 ## DATA OUTPUT PREFERENCES ## 1. Prioritize TEXT format (print) for statistical information. 991 2. Print a description before outputting statistics. 992 3. Generate images only if text description is inadequate. 993 ## METHODOLOGY REQUIREMENTS ## 994 1. Provide highly detailed methods, especially for data cleaning. 2. Specify actions for each feature without omissions. 995 996 ## RESOURCE MANAGEMENT ## 1. Consider runtime and efficiency, particularly for: 997 - Data visualization 998 - Large dataset handling - Complex algorithms 999 2. Limit generated images to a MAXIMUM of 10 for EDA. 1000 3. Focus on critical visualizations with valuable insights. 1001 ## OPTIMIZATION EXAMPLE ## 1002 When using seaborn or matplotlib for large datasets: - Turn off unnecessary details (e.g., annot=False in heatmaps) 1003 - Prioritize efficiency in plot generation 1004 ############# 1005 # TASK # {task} 1007 ############# 1008 # RESPONSE # Let's work this out in a step by step way. 1009 1010 ############## # START PLANNING # 1011 Before you begin, please request the following documents from me, which contain 1012 important information that will guide your planning: 1. Report and plan from the previous phase 1013 2. Available tools in this phase 1014 3. Sample data for analysis 1015 1016 Please design plan that is clear and specific to each FEATURE for the current development phase: {phase_name}. 1017 The developer will execute tasks based on your plan. 1018 I will provide you with INFORMATION, RESOURCE CONSTRAINTS, and previous reports and plans. 1019 You can use the following reasoning pattern to design the plan: 1020 1. Break down the task into smaller steps. 2. For each step, ask yourself and answer: 1021 - "What is the objective of this step?" 1022 - "What are the essential actions to achieve the objective?" - "What features are involved in each action?" 1023 - "Which tool can be used for each action? What are the parameters of the tool?" 1024 - "What are the expected output of each action? What is the impact of the action on the data?" 1025

Prompt of Agent Developer / Task Prompt

1026 D.1.4 DEVELOPER

tests.

1028Developer is responsible for implementing and debugging code based on the structured plans1029generated by the Planner. The Developer's key function is to translate the high-level task roadmap1030into executable code, resolve any arising issues, and perform unit tests to ensure the functionality of1031the solution.

Description: The Developer first reviews the task plan and the relevant competition information. It can gathers

code from previous phases when necessary and uses LLMs

to generate new code. The Developer also cleans up any

redundant code sections, writes functions, and ensures the

code runs correctly by debugging and performing unit tests.

It iterates through the process until the code passes all

Role: write and implement code according to plan

20

Prompt of Agent Developer / Task Prompt
The off the rest of the rest o
<pre># CONTEXT # {phases in context}</pre>
Currently, I am at phase: {phase_name}.

INFORMATION
{background_info}
{state_info}

PLAN
{plan}

TASK
{LdSK}
RESPONSE: BLOCK (CODE & EXPLANATION) # TASK 1:
THOUGHT PROCESS:
[Explain your approach and reasoning] CODE:
" python
EXPLANATION:
[Brief explanation of the code and its purpose]
TASK 2:
[Repeat the above structure for each task/subtask]

START CODING
Before you begin, please request the following information from me:
2. All features of the data
3. Available tools
Once you have this information, provide your complete response with code and
explanations for all tasks in a single message.
Develop an efficient solution based on the Planner's provided plan:
 Implement specific tasks and methods outlined in the plan Ensure code is clear, concise, and well-documented
3. Utilize available tools by calling them with correct parameters
 Consider data types, project requirements, and resource constraints Write code that is easily understandable by others
Remember to balance efficiency with readability and maintainability.

D.1.5 REVIEWER

Reviewer is responsible for evaluating the performance of other agents in completing tasks and providing constructive feedback.

Prompt of Agent Reviewer / Task Prompt

Role: assess agent performance and offer feedback Description: The Reviewer agent evaluates the performance of multiple agents. In each evaluation phase, it merges suggestions and scores from different agents into a unified report. It interacts with a LLM to generate detailed feedback, iterating through rounds to assess task results, merging agent responses, and producing both final scores and constructive suggestions.

Prompt of Agent Reviewer

CONTEXT

{phases_in_context}

############

TASK

Tour cask is to assess the periormance of several agenes in compreting mase. (
phase_name}.
I will provide descriptions of each agent, the tasks they performed, and the outcomes
of those tasks.
Please assign a score from 1 to 5 for each agent, with 1 indicating very poor
performance and 5 indicating excellent performance.
Additionally, provide specific suggestions for improving each agent's performance, if applicable.
If an agent's performance is satisfactory, no suggestions are necessary.

Each phase involves collaboration between multiple agents. You are currently

evaluating the performance of agents in Phase: {phase_name}.

RESPONSE: JSON FORMAT # Let's work this out in a step by step way.

D.1.6 SUMMARIZER

Summarizer is responsible for generating summaries, designing questions, and reorganizing both questions and answers to produce structured reports based on the competition phases.

Prompt of Agent Summarizer / Task Prompt

Role: assess agent performance and offer feedback Description: The agent Summarizer works through various phases, each focusing on a specific task like choosing relevant images, designing key questions, answering phase-related questions, and organizing the responses into a structured report. Each phase involves interaction with provided inputs such as competition information, the planner's plan, and the reviewer's evaluation to synthesize the most relevant insights.

	# TASK #
	Please reorganize the answers that you have given in the previous step, and synt
	tnem into a report.
	# RESPONSE: MARKDOWN FORMAT #
	# REPORT
	## QUESTIONS AND ANSWERS
	### Question 1 What files did you process? Which files were generated? Answer with detailed fil
	### Answer 1
	[answer to question 1]
	### Question 2
	Which features were involved in this phase? What changes did they undergo? If an feature types were modified answer which features are modified and how the
	modified. If any features were deleted or created, answer which features are
	deleted or created and provide detailed explanations. (This is a FIXED ques
	### Answer 2
	[answer to question 2]
	### Question 3
	[repeat question 3]
	### Answer 3
	Tauswer to Anestron ol
	### Question 4
	[repeat question 4] ### Answer 4
	[answer to question 4]
	### Ouestion 5
	[repeat question 5]
	### Answer 5
	[answer to question 5]
	### Question 6
	[repeat question 6]
	[answer to question 6]
ļ	****
	# START REORGANIZE OUESTIONS #

D.2 UNIT TESTS

In data science competitions, code generated by agents may be executable in the Python interpreter, but this execution does not guarantee correctness. To ensure that data dependencies are properly handled, a Unit Test Tool is necessary. In our research, where the framework operates iteratively, we aim to separate tasks corresponding to different states in data science competitions. Each phase builds upon the results of the previous one, making it crucial to confirm that logic remains sound, data processing is accurate, and information transfers seamlessly from one state to the next. Our Unit Test Tool plays a key role in supporting the self-refine phase of LLM agents.

We developed unit tests (in the accompanying Table 6) based on issues identified during the execution of weak baseline, strong baseline and our AutoKaggle. If the code fails to run in the Python interpreter, an error message is relayed to the agent Reviewer. If the code passes this initial stage, it progresses to the Unit Test Tool, where all required tests are executed in a loop. If a test fails, the reason is logged as short-term memory and passed to the next review state. The review and planning stages work in an adversarial interaction: the review phase compiles the reasons for failed unit tests, while the planner addresses these failures in subsequent iterations.

1241

Table 6: Overview of unit tests for state DC, FE, and MBVP. These unit tests handle to detect missing values, outliers, duplicates, and other data consistency issues.

State	Unit test name	Unit test description
	test_document_exist	Test if cleaned_train.csv and cleaned_test.csv data exist.
State DC	test_no_duplicate_cleaned_train	Test if there are any duplicate rows in the cleaned_train.csv.
State DC	test_no_duplicate_cleaned_test	Test if there are any duplicate rows in the cleaned_test.csv.
	test_readable_cleaned_train	Test if the cleaned_train.csv is readable.
	test_readable_cleaned_test	Test if the cleaned_test.csv is readable.
	test_cleaned_train_no_missing_values	Test if the cleaned_train.csv contains missing value.
	test_cleaned_test_no_missing_values	Test if the cleaned_test.csv contains missing value.
	test_cleaned_train_no_duplicated_features	Test if the cleaned_train.csv contains duplicate features.
	test_cleaned_test_no_duplicated_features	Test if the cleaned_test.csv contains duplicate features.
	test_cleaned_difference_train_test_columns	Test if the cleaned_train.csv and cleaned_test.csv have the same features except for target variable.
	test_cleaned_train_no_missing_target	Test if the target variable is in cleaned_train.csv.
	test_document_exist	Test if processed_train.csv and pro- cessed test csv data exist
State FE	test_processed_train_feature_number	Test if the feature engineering phase is per- formed well in processed_train.csv.
	test_processed_test_feature_number	Test if the feature engineering phase is per- formed well in processed_test.csv.
	test_file_size	Test if processed data is larger than a threshold.
	test_processed_train_no_duplicated_features	Test if the processed_train.csv contains dupli- cate features.
	test_processed_test_no_duplicated_features	Test if the processed_test.csv contains duplicate features.
	test_processed_difference_train_test_columns	Test if the processed_train.csv and pro- cessed_test.csv have the same features except for target varibale.
	test_processed_train_no_missing_target	Test if the target variable is in pro- cessed_train.csv.
	test_document_exist	Test if a submission file exists.
	test_no_duplicate_submission	Test if there are any duplicate rows in the sub- mission file
	test_readable_submission	test if the submission file is readable.
State MRVP	test_file_num_submission	Test if the submission file and sam-
		ple_submission.csv have the same number of rows.
	test_column_names_submission	Test if the submission file and sam- ple_submission.csv have the same column names.
	test_submission_validity	 Test if the submission file and sample_submission.csv have the same data index. 2) Test if the submission file and sample_submission.csv have the same numerical range

1296 D.3 MACHINE LEARNING TOOLS DETAILS

1297

Table 7: Overview of Tools for state DC, FE, and MBVP. This table presents various tools categorized by their functionality.

State	Tool name	Tool description
	FillMissingValues	Fills missing values or removes columns with missing values based on a threshold.
	RemoveColumns WithMissingData	Removes columns containing missing values from DataFrame based on a threshold.
State DC	DetectAndHandleOutliersZscore	Detects and handles outliers in specified columns us the Z-score method.
State DC	DetectAndHandleOutliersIqr	Detects and handles outliers in specified columns us the Interquartile Range (IQR) method.
	RemoveDuplicates	Removes duplicate rows from a DataFrame.
	ConvertDataTypes	Converts the data type of specified columns in DataFrame.
	FormatDatetime	Formats datetime columns to a specified format.
		Performs one-hot encoding on specified categor
	OneHotEncode	columns.
	LabelEncode	Performs label encoding on specified categor columns.
	FrequencyEncode	Performs frequency encoding on specified categor columns.
	TargetEncode	Performs target encoding on specified categor columns.
	CorrelationFeatureSelection	Performs feature selection based on correlation and sis.
State FE	VarianceFeatureSelection	Performs feature selection based on variance analys
	ScaleFeatures	Scales numerical features in the specified columns DataFrame.
	PerformPca	Performs Principal Component Analysis (PCA) on specified columns of a DataFrame.
	PerformRfe	Performs Recursive Feature Elimination (RFE) on specified columns of a DataFrame.
	CreatePolynomialFeatures	Creates polynomial features from specified column a DataFrame.
	CreateFeatureCombinations	Creates feature combinations from specified column a DataFrame.
State MBVP	TrainAndValidation AndSelectTheBestModel	Trains, evaluates, and selects the best machine learn model based on the training data and labels, return the best performing model along with the performan scores of each model and their best hyperparameters

Examples of Tool Schema. In this paper, we provide two schema formats for each machine learning tool: JSON and Markdown. Here, we take the FillMissingValues tool as an example and provide schemas in both formats.

Description: Fill missing values in specified columns of a DataFrame. This tool can h	
Description: Fill missing values in specified columns of a DataFrame. This tool can h	11
both numerical and categorical features by using different filling methods	andle
1357 Applicable Situations: Handle missing values in various types of features	
1358 Parameters:	
1359 • data:	
1360 Gata.	
1361 – Type: pd. DataFrame	
- Description: A pandas DataFrame object representing the dataset.	
1363 • columns:	
- Type: string array	
Description: The name(s) of the column(s) where missing values show	ld be
1367 filled.	
1368 • method:	
1369 – Type: string	
1370 – Description: The method to use for filling missing values	
1371 – Enum: auto mean median mode constant	
$1372 \qquad - Default auto$	
1373	
• IIII_value:	
1375 - Type: number string null	
1376 – Description: The value to use when method is constant.	
1377 – Default: None	
1378 Required: data, columns	
Result: Successfully fill missing values in the specified column(s) of data.	
1380 Notes:	
• The auto method uses mean for numeric columns and mode for non-num	neric
1382 columns.	
• Using mean or median on non-numeric columns will raise an error.	
• The mode method uses the most frequent value, which may not always be a	ppro-
priate.	
• Filling missing values can introduce bias, especially if the data is not missing	com-
1388 pletely at random.	
• Consider the impact of filling missing values on your analysis and model p	erfor-
1390 mance.	
1391	
1392	
1393	
1394	
1395	
1396	
1397	
1398	
1400	
1400	
1402	
1403	

```
JSON-formatted tool schema for FillMissingValues
1405
1406
         {
1407
             "name": "fill_missing_values",
1408
             "description": "Fill missing values in specified columns
1409
            of a DataFrame. This tool can handle both numerical and
1410
            categorical features by using different filling methods
1411
             .",
1412
             "applicable_situations": "handle missing values in
1413
            various types of features",
             "parameters": {
1414
                  "data": {
1415
                      "type": "pd.DataFrame",
1416
                      "description": "A pandas DataFrame object
1417
             representing the dataset."
1418
                  },
1419
                  "columns": {
1420
                      "type": ["string", "array"],
1421
                      "items": {
1422
                          "type": "string"
1423
                      },
                      "description": "The name(s) of the column(s)
1424
            where missing values should be filled."
1425
                  },
1426
                  "method": {
1427
                      "type": "string",
1428
                      "description": "The method to use for filling
1429
            missing values.",
1430
                      "enum": ["auto", "mean", "median", "mode", "
1431
             constant"],
1432
                      "default": "auto"
1433
                  },
1434
                  "fill_value": {
                      "type": ["number", "string", "null"],
1435
                      "description": "The value to use when method is '
1436
             constant'.",
1437
                      "default": null
1438
                  }
1439
             },
1440
             "required": ["data", "columns"],
1441
             "result": "Successfully fill missing values in the
1442
             specified column(s) of data",
1443
             "additionalProperties": false,
1444
             "notes": [
1445
                  "The 'auto' method uses mean for numeric columns and
            mode for non-numeric columns.",
1446
                  "Using 'mean' or 'median' on non-numeric columns will
1447
              raise an error.",
1448
                  "The 'mode' method uses the most frequent value,
1449
             which may not always be appropriate.",
1450
                  "Filling missing values can introduce bias,
1451
            especially if the data is not missing completely at
1452
            random.",
1453
                  "Consider the impact of filling missing values on
1454
            your analysis and model performance."
1455
             ]
1456
1457
```

Tool use. During execution, we extract the machine learning tools specified in the plan generated by Planner and use them as queries to search the entire documentation of machine learning tools. Since the plan includes multiple tools, we retrieve several tools based on their similarity to the queries. The Developer then uses the retrieved tools to carry out the task.

1463 D.4 TOOL UTILIZATION

In the multi-agent framework designed for autonomous data science tasks, tools serve not only as automation resources but also as integral components of the workflow. The framework enables agents to dynamically access and execute tools as they transition through various problem-solving states, ensuring adaptability and efficiency.

1469 The tool utilization process in this framework is structured around a systematic approach. Tool in-1470 formation is first stored in the system's Memory, which is implemented as a vector database. This Memory holds detailed explanations regarding each tool's functionality, usage, and context. A con-1471 figuration file is used to map specific tools to the states in which they are required, allowing agents 1472 to reference and identify the appropriate tools at each stage of the problem-solving process. To de-1473 termine which tools are required in each state, the table 7 provides an overview of tools categorized 1474 by their functionality. As an agent moves into a particular state, it consults the configuration file to 1475 determine the relevant tools. From the figure 1 shown, the agent subsequently queries the Memory 1476 to retrieve detailed explanations for the tool's use, and finally, executes the tool with precision based 1477 on the retrieved information. 1478

This dynamic interaction between the Memory, configuration file, and agents facilitates seamless tool integration, empowering agents to operate autonomously while maintaining flexibility and ensuring accurate tool application throughout the autonomous process.

1482

1462

1464

1483 D.5 USER INTERACTION

1484

At each stage of problem-solving, two Human-in-the-Loop methods are employed. Before the Planner formulates a plan, human can interact with the command line. The input consists of meticulously manually crafted rules, each one carefully cataloged in a handbook. Memory module subsequently retrieved these predefined rules, integrating this human-driven knowledge in prompt engineering to guide the Planner's next steps. After generating the plan, humans can review and and refine the Planner's output. They inspect areas where the logical flow seems inconsistent, focusing particularly on points where the output diverges from reality to address hallucination issues.

1491

¹⁴⁹² E CASE STUDY: TITANIC

1494 E.1 BACKGROUND UNDERSTANDING 1495

In this step, the system employs a LLM (GPT-40) to extract and summarize the key information from
 the Titanic Kaggle competition. Upon completion of this process, a markdown file is automatically
 generated containing essential competition details, which include the competition overview, dataset
 information, and evaluation metrics. Below is an excerpt of the generated markdown file:

- 1500
- 1501 1502
- 1503
- 1504
- 1505
- 1506
- 1502
- 1508 1509
- 1510
- 1511

Background Understanding

Competition Information

1. Competition Overview

ground truth).

3. Problem Definition

6. Evaluation Metrics

7. Submission Format

8. Other Key Aspects

the model.

containing:

need to be predicted.

2. Files



1536 1537

1538 1539 1540

1541

E.2 PRELIMINARY EXPLORATORY DATA ANALYSIS

'PassengerId': Sorted in any order.

1542 In this state, an autonomous exploratory analysis is conducted to understand the Titanic dataset. The 1543 LLM will plan a process: 1) the dataset is loaded and inspected to reveal its structure, including data types, basic statistics, and missing values. A univariate analysis follows, where the distributions 1544 of key numerical features like age and fare are explored, as well as the frequencies of categorical 1545 features such as passenger class and gender. 2) A bivariate analysis is performed to examine the 1546 relationship between each feature and the target variable, survival. This step involves visualizing 1547 the survival rates across different categories and examining how numerical features vary between 1548 survivors and non-survivors. 3) A correlation analysis investigates relationships between numerical 1549 features, visualized through a correlation matrix and additional pairwise comparisons. This analysis 1550 provides a deeper understanding of feature interactions, highlighting those that may have the most 1551 significant impact on survival. Below is an excerpt of the generated file:

The Titanic competition challenges participants to build a predictive model to

determine which passengers survived the shipwreck of the RMS Titanic on April 15,

1912, using passenger data such as name, age, gender, and socio-economic class.

train.csv: Contains data for 891 passengers, including whether they survived (

test.csv: Includes data for 418 passengers, but without survival outcomes-these

sample_submission.csv: A template showing the required format for submissions.

The objective is to predict the survival of passengers based on various features in the dataset. Participants build a model using the training set to learn from the

The evaluation metric is **accuracy**: the percentage of correct predictions made by

- It's beneficial to explore shared notebooks for additional insights and approaches.

Submissions must be in CSV format with exactly 418 entries plus a header row,

- 'Survived': Binary predictions (1 for survived, 0 for deceased).

The competition runs indefinitely with a rolling leaderboard.

- Kaggle provides a community forum for discussion and collaboration.

data and apply this model to make predictions on the test set.

- 1552
- 1553
- 1554 1555
- 1556
- 1557
- 1558

1559 1560

- 1560
- 1562
- 1563

1564

,	Great thank you for the information. Let's move forward with a detailed plan for
	Preliminary Exploratory Data Analysis phase. The main objectives of this pha
	to understand the data structure, examine distributions, and identify any i
	patterns or anomalies. We will break down the task into four key steps:
-	### 1. Understanding Data Structure
-	#### Objective:
	To get a high-level overview of the data types, missing values, and basic statis
-	#### Actions:
-	<pre>1. **Load the Data**: - Use Pandas to read the `train.csy` and `test.csy` files.</pre>
	- Tools: 'pd.read_csv()'
,	- Expected Output: DataFrames for train and test data.
í	
	<pre>3. **Basic Statistical Summary**:</pre>
-	### 2. Univariate Analysis of Numerical Features
1	#### Jecclve: To understand the distribution of each numerical feature.
Ĵ	
-	#### Actions: 1. **Histograms for Numerical Features**:
Î	·····
-	2. **Box Plots for Numerical Features**:
-	### 3. Analysis of Categorical Features
1	#### Objective: To examine the distribution and frequency of each categorical feature.
-	#### Actions: 1. **Bar Plots for Categorical Features***
	·····
-	2. **Value Counts**:
-	#### 4. Initial Insights and Summary
	To summarize the initial findings and identify any immediate data quality issues
	patterns.
	#### Actions:
-	1. **Summarize Key Findings**:
	2. **Identify Next Steps**:
	### Constraints:
-	- **Efficiency**: Ensure that visualizations and computations are efficient and
	reproducible.
	understanding.
	### Exported Outputs.
	1. DataFrames for train and test sets.
	2. Summary of data types and missing values.
	э. вазие statistical summaries. 4. Histograms and box plots for numerical features.
	5. Bar plots and frequency distributions for categorical features.
4	6. A detailed summary of initial insights and recommended next steps.
F	By following this plan, you'll gain a comprehensive preliminary understanding of
	data, setting a solid foundation for subsequent phases.

1620 Preliminary Exploratory Data Analysis Code 1621 1622 import os 1623 import pandas as pd 1624 import matplotlib.pyplot as plt import seaborn as sns 1625 1626 # File paths train_file_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/ 1627 titanic/train.csv' 1628 test_file_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/ titanic/test.csv' 1629 1630 # Loading datasets train_df = pd.read_csv(train_file_path) 1631 test_df = pd.read_csv(test_file_path) 1632 # Output dataset shapes 1633 print(f"Train dataset shape: {train_df.shape}") 1634 print(f"Test dataset shape: {test_df.shape}") 1635 # Inspect data types and missing values 1636 print("Missing values in train dataset:") print(train_df.isnull().sum()) 1637 print("\nMissing values in test dataset:") 1638 print(test_df.isnull().sum()) 1639 # Basic statistical summary for numerical features 1640 print("Statistical summary for numerical features in train dataset:") print(train_df.describe()) 1641 1642 # Plot histograms for numerical features numerical_features = ['Age', 'Fare', 'SibSp', 'Parch'] 1643 for feature in numerical_features: 1644 plt.figure(figsize=(8, 4)) sns.histplot(train_df[feature], bins=30, kde=True) 1645 plt.title(f'Histogram of {feature}') 1646 plt.xlabel(feature) plt.ylabel('Frequency') 1647 plt.savefig(f'/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/ 1648 titanic/pre_eda/images/histogram_{feature}.png') plt.close() 1649 1650 # Plot bar charts for categorical features categorical_features = ['Pclass', 'Sex', 'Embarked', 'Survived'] 1651 for feature in categorical_features: 1652 plt.figure(figsize=(8, 4)) sns.countplot(x=train_df[feature]) 1653 plt.title(f'Bar chart of {feature}') 1654 plt.xlabel(feature) plt.ylabel('Count') 1655 plt.savefig(f'/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/ 1656 titanic/pre_eda/images/bar_chart_{feature}.png') plt.close() 1657 1658 # Summarize key findings key_findings = """ 1659 Key Findings from Preliminary EDA: 1660 1. Missing Values: - Age has 177 missing values in the train dataset. 1661 - Cabin has significant missing values. 1662 - Embarked has 2 missing values in the train dataset. 1663 2. Numerical Features: 1664 - Age distribution shows a right skew. - Fare distribution is highly skewed. 1665 1666 3. Categorical Features: - More males than females in the train dataset. 1667 - Most passengers embarked from port $^{\prime}\,\text{S}^{\prime}$ (Southampton). 1668 - Survival rate is around 38% in the train dataset. 1669 print(key_findings) 1670 1671

Preliminary EDA

Train Dataset Hea	ad:				
PassengerId	Survived	Pclass	. Fare	Cabin	Embarke
0 I 1 2	0	3	71 2833	C85	S
2 3	1	3	7.9250	NaN	S
3 4	1	1	53.1000	C123	S
4 5	0	3	8.0500	NaN	S
[5 rows x 12 colu	umns]				
Test Dataset Head	d:				
PassengerId	Pclass	. Cabin Em	barked		
0 892	3	NaN	Q		
2 893	3	NaN	S		
2 895 3 895	3	NaN	S		
4 896	3	NaN	S		
[5 rows x 11 col)	um ne l				
[5 10w3 x 11 co10					
Missing Values in	n Test Data	set:			
PassengerId	0				
Pclass	0				
Sex	0				
Age 8	36				
SibSp	0				
Parch	0				
Fare	1				
Cabin 32	27				
Embarked	0				
dtype: int64					
Frequency of Pcla	ass:				
Pclass 3 491					
1 216					
Z 184 Name: count, dtyp	pe: int64				
Frequency of Cov					
Sex	•				
male 577					
remale 314 Name: count, dtyp	pe: int64				
Frequency of Emba Embarked S 644 C 168 Q 77	arked:				

E.3 DATA CLEANING

We demonstrate the data analysis capabilities of our framework using the age column from the Titanic competition's training set as an example. In the pre-EDA phase, the distribution of the age histogram is as shown in Figure 6. During the data cleaning phase, we filter out missing values using unit tests. You can see a comparison of the age box plots before and after the outliers have been processed in Figure 7. In the deep-EDA phase, the distribution of the age histogram is as shown in Figure 8.



Figure 7: Comparison of age box plots before and after treatment of outliers. The image on the left is before the outliers are processed and the one on the right is after the process is done



Figure 8: The histogram of age after outliers are processed

Data Cleaning Code

```
1783
1784
            import sys
1785
            import os
1786
            import numpy as np
            import pandas as pd
1787
            import matplotlib.pyplot as plt
1788
            import seaborn as sns
1789
            1790
            sys.path.append(os.path.dirname(os.path.abspath(__file__)))
1791
            from tools.ml_tools import *
1792
            def generated_code_function():
1793
                # Load datasets
1794
                train_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/
                titanic/train.csv'
1795
               test_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/
1796
                titanic/test.csv'
                train_df = pd.read_csv(train_path)
1797
                test_df = pd.read_csv(test_path)
1798
                # Handle missing values
1799
               for df in [train_df, test_df]:
1800
                   df = fill_missing_values(df, columns=['Age', 'Fare'], method='median')
                   df = fill_missing_values(df, columns=['Embarked'], method='mode')
1801
                   df = remove_columns_with_missing_data(df, columns=['Cabin'])
1802
                # Convert data types
1803
                for df in [train_df, test_df]:
1804
                   df = convert_data_types(df, columns=['PassengerId', 'Pclass'], target_type='
                str')
1805
               train_df = convert_data_types(train_df, columns=['Survived'], target_type='str')
1806
                # Plot outliers and handle using IQR method
1807
               def plot_outliers(data, columns, suffix):
1808
                   output_dir = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/
                titanic/data_cleaning/images/'
1809
                   os.makedirs(output dir, exist ok=True)
1810
                   for column in columns:
                       plt.figure(figsize=(10, 5))
1811
                       sns.boxplot(x=data[column])
1812
                       plt.title(f'Boxplot of {column} {suffix}')
                       plt.savefig(f'{output_dir}{column}_{suffix}.png')
1813
                       plt.close()
1814
               columns_with_outliers = ['Age', 'Fare']
1815
               plot_outliers(train_df, columns_with_outliers, 'before_outliers')
1816
               for df in [train_df, test_df]:
1817
                   df = detect_and_handle_outliers_iqr(df, columns=columns_with_outliers, factor
1818
                =1.5, method='clip')
1819
               plot_outliers(train_df, columns_with_outliers, 'after_outliers')
1820
                # Save cleaned datasets
1821
               train_df.to_csv('/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/
1822
                titanic/cleaned_train.csv', index=False)
                test_df.to_csv('/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/
1823
                titanic/cleaned_test.csv', index=False)
1824
                _name__ == "__main__":
            if
1825
               generated_code_function()
1826
1827
1828
1829
1830
```

- 1831
- 1832
- 1833
- 1834
- 1835

1836	
1000	
	Data Cleaning Result
1007	Data Creaning restart

<pre> Sex 0 Age 177 Cabin 687 Embarked 2 dtype: int64 Wissing values in test dataset before handling: PassengerId 0 Age 66 Fare 1 Cabin 327 Embarked 0 dtype: int64 Wissing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Wissing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Parch int64 PassengerId object Sex 0bject Starp float64 Embarked 0 Survived object Ticket 0 dtype: object Data types in test dataset after conversion: Age float64 Embarked 0 Survived object Parch int64 PassengerId object Parch int64 PassengerId object Sex 0bject Starp in test dataset after conversion: Age float64 Embarked object Parch int64 PassengerId o</pre>	Missing values PassengerId	in train dataset before handling: 0
Jek 0 Age 177 Cabin 687 Embarked 2 dtype: int64 Missing values in test dataset before handling: Passenger14 0 Age 8 Tare 1 Cabin 327 Embarked 0 Comparison of the set of		0
<pre>''''''''''''''''''''''''''''''''''''</pre>	Sex	U 177
Cabin 667 Embarked 2 dtype: int64 Missing values in test dataset before handling: PassengerId 0 Age 86 Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 PassengerId object FassengerId object Six object Six object Data types in test dataset after conversion: Age float64 Embarked object FassengerId object FassengerId object Six object Data types in test dataset after conversion: Age float64 Embarked object FassengerId object FassengerId object Six object Data types in test dataset after conversion: Age float64 Embarked object FassengerId object Cleass object Data types in test dataset after conversion: Age float64 Embarked object FassengerId object Cleass object Six object Six object Cleass object Six object Cleas object Cleased tast asved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competiton/titenic/cleaned_test.cxv		111
<pre>Embarked 2 dtype: int64 Missing values in test dataset before handling: PassengerId 0 Age 86 Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SiNSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SiNSp 0 Ticket 0 dtype: int64 Missing values in test dataset after conversion: Age float64 Missing float64 Mame 00ject Parch 1nt64 PassengerId 00ject Sex 00ject Sex 00ject Sex 00ject Lata types in test dataset after conversion: Age float64 Mame 00ject Parch 1nt64 PassengerId 00ject Parl 1nt64 Survived 00ject Sex 00ject Sex 00ject Sex 00ject Sex 00ject Cabin Sing 1nt64 Survived 00ject Parl 1nt64 PassengerId 00ject Parl 1nt64 PassengerId 00ject Cabin Sing 1nt64 Survived 00ject Cabin Sing 1nt64 PassengerId 00ject Cabin Sing 1nt64 PassengerId 00ject Cabin Sing 1nt64 Survived 00ject Cabin Sing 1nt64 PassengerId 00ject Cabin Sing 1nt64 PassengerId</pre>	Cabin	687
dtype: int64 Missing values in test dataset before handling: PassengerId 0 Age 86 Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 PassengerId object PassengerId object SibSp 1 Ticket 0 dtype: object At types in test dataset after conversion: Age float64 Embarked object Fare float64 Embarked object Fare float64 Mame object Stat uppes in test dataset after conversion: Age float64 Embarked object Fare float64 Mame object Ticket 0 dtype: object Cleaned tast as aved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competiton/ticmic/cleaned_test.csv	Embarked	2
Missing values in test dataset before handling: PassengerId 0 Age 86 Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SiSSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SiSSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 PassengerId object Pclass object SixSp int64 Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SixSp int64 Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SixSp int64 Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SixSp int64 Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SixSp int64 PassengerId object Pclass object SixSp int64 PassengerId object Parch int64 Parc	dtype: int64	
<pre>PassengerId 0 Age 86 Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Pats doject SibSp int64 Data types in test dataset after conversion: Age float64 Name object Pats doject Ticket 0 dtype: object Data types in test dataset after conversion: Age float64 Survived object Ticket 0 Data types in test dataset after conversion: Age float64 Survived object Fare float64 Survived object Ticket 0 Data types in test dataset after conversion: Age float64 Name object Pclass object SibSp int64 Survived object Ticket 0 Data types in test dataset after conversion: Age float64 Name object Pclass object SibSp int64 Survived object Ticket 0 Data types in test dataset after conversion: Age float64 Name object Pclass object Cleaned taining data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competition/tlainc/oleaned_test.cov Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competition/tlainc/oleaned_test.cov Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competition/tlainc/oleaned_test.cov </pre>	Missing values	in test dataset before handling:
<pre> Age 86 Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SlSSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SlSSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Parch int64 PassengerId object Pclass object SlSSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SlSSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SlSSp int64 Survived object Ticket object Pclass object Parch int64 PassengerI object Parch int64 PassengeI object Parch int64 PassengeI object Pclass object</pre>	PassengerId	0
AgesbFare1Cabin327Embarked0dtype: int640Missing values in train dataset after handling:Age0Embarked0SibSp0Ticket0Missing values in test dataset after handling:Age0Embarked0		
<pre>Fare 1 Cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Pacha float64 Survived object Data types in test dataset after conversion: Age float64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Survived object Cleaned test dataset after conversion: Age float64 Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titaic/cleaned_train.csv Cleaned test data saved to /mn</pre>	Age	80
<pre>cabin 327 Embarked 0 dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SiNSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SiNSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 ParsengerId object Pareh int64 PassengerId object SiNSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked colject Pareh int64 PassengerId object Ticket object Data types in test dataset after conversion: Age float64 Survived object Pclass object SiNSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Pclass object Sex object Sex object Sex object Class object Sex object Sex object Class object Sex object Class object Class object SiNSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv</pre>	Fare	1
Embarked0dtype: int64Missing values in train dataset after handling:Age0Embarked0SibSp0Ticket0dtype: int64Missing values in test dataset after handling:Age0Embarked0SibSp0Ticket0dtype: int64Data types in train dataset after conversion:Agefloat64EmbarkedobjectParenfloat64Parentint64PassengerIdobjectPclassobjectSibSpint64SurvivedobjectTicketobjectTicketobjectParentint64PasengerIdobjectTicketobjectTicketobjectParentint64PasengerIdobjectParentint64PasengerIdobjectParentint64PasengerIdobjectPassengerIdobjectStasobjectStasobjectStasobjectStasobjectParentint64PasengerIdobjectPasengerIdobjectStasobjectStasobjectStasobjectStasobjectCleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_age	Cabin	327
dtype: int64 Missing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object SibSp int64 Survived object Ticket object Ticket object Data types in test dataset after conversion: Age float64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Ticket object Ticket object Ticket object Parch int64 PassengerId object Parch int64 PassengerId object Ticket object Ticket object SiSSp int64 Ticket object SiSSp int64 Ticket object SiSSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agent competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Embarked	0
Missing values in train dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 PassengerId object Pclass object Skx object SibSp int64 Survived object Ticket object Ticket object Ticket object Data types in test dataset after conversion: Age float64 PassengerId object Pclass object SibSp int64 Survived object Ticket object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Ticket object Ticket object Ticket object Pclass object SibSp int64 Survived object Ticket object Ticket object Pclass object Pclass object Pclass object Pclass object Ticket object Pclass	dtype: int64	
Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 0 Missing values in test dataset after handling: Age Age 0 Embarked 0 SibSp SibSp 0 Ticket 0 dtype: int64 0 Data types in train dataset after conversion: Age float64 Pata types in train dataset after conversion: Age float64 Name object Parch int64 PassengerId object Pclass object Ticket object Vivived object Ticket object Ticket object Parch int64 PassengerId object Parce float64 Mame object Parce float64 Name object Parce float64 RasengerId object	Missing values	in train dataset after handling:
Embarked0SibSp0Ticket0dtype: int64Missing values in test dataset after handling: Age0Embarked0Embarked0SibSp0Ticket0dtype: int64Data types in train dataset after conversion: Agefloat64RmbarkedobjectFarefloat64NameobjectParchint64PassengerIdobjectSibSpint64SurvivedobjectTicketobjectPata types in test dataset after conversion: AgeAgefloat64NameobjectPata types in test dataset after conversion: AgeAgefloat64EmbarkedobjectTicketobjectParchint64SurvivedobjectPata types in test dataset after conversion: AgeAgefloat64EmbarkedobjectParchint64PassengerIdobjectParasobjectStSpint64PassengerIdobjectStSpint64PassengerIdobjectParasobjectParasobjectCleaned taining data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competition/titanic/cleaned_train.csvCleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Age	0
<pre>>: SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Pare float64 Name object Pares object SibSp int64 Survived object Data types in test dataset after conversion: Age float64 Embarked object Pares object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Embarked object Data types in test dataset after conversion: Age float64 Survived object Data types in test dataset after conversion: Age float64 Embarked object Data types in test dataset after conversion: Age float64 Embarked object Data types in test dataset after conversion: Age float64 Embarked object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv </pre>	Embarked	0
<pre> SibSp 0 Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Pambarked object Fare float64 Name object Pclass object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Survived object Pclass object Data types in test dataset after conversion: Age float64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Survived object Class object Class object SibSp int64 Survived object Tare float64 Name object Pclass object Classe object Classe object SibSp int64 Ticket object Claened training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv </pre>		
Ticket 0 dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Pare float64 Embarked object Pare float64 Embarked object Fare float64 Survived object Tare float64 Survived object Pare float64 Embarked object Fare float64 Cases object SibSp int64 SibSp int64 Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	sibsp	0
dtype: int64 Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Parch int64 PassengerId object Parch int64 PassengerId object Fare float64 Embarked object Fare float64 PassengerId object Parch int64 PassengerId object Parch int64 PassengerId object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Ticket	0
Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Pare float64 Competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	dtype: int64	
Missing values in test dataset after handling: Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Embarked object Fare float64 Embarked object Parch int64 PassengerId object Fare float64 Embarked object Fare float64 Embarked object Fare float64 Embarked object Parsh int64 PassengerId object Parsh int64 PassengerId object Parsh int64 Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv		
Age 0 Embarked 0 SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Parch int64 PassengerId object Parch int64 PassengerId object Parch int64 PassengerId object Parch int64 PassengerId object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Missing values	in test dataset after handling:
<pre>SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pare float64 PassengerId object Pclass object SibSp int64 Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv</pre>	Embarked	0
<pre>SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object SibSp int64 Survived object Ticket object dtype: object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Parch int64 PassengerId object Pclass object SibSp int64 SibSp int64 Claase object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv</pre>		•
SibSp 0 Ticket 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Embarked object Fare float64 Parch int64 PassengerId object Parch int64 PassengerId object Sex object Sex object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv		
<pre>11Ckt 0 dtype: int64 Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object SibSp int64 Survived object Class object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Embarked object Fare float64 Embarked object Parch int64 PassengerId object Parch int64 PassengerId object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv</pre>	SibSp	0
Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object Sex object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	dtype: int64	U
Data types in train dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Embarked object Parch int64 PassengerId object Parch int64 PassengerId object Sex object Sex object Sex object Sex object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	acipo. Incor	
Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object dtype: object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Parch int64 PassengerId object Sex object Sex object Sex object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Data types in '	train dataset after conversion:
Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object SibSp int64 Survived object Ticket object dtype: object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Sex object Sex object Sex object Sex object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Age Data a via a d	float64
Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Parch float64 PassengerId object Parch int64 PassengerId object Pclass object Sex object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Embarked	object float64
Parch int64 PassengerId object Pclass object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Parch float64 PassengerId object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Name	object
PassengerId object Pclass object Sex object SibSp int64 Survived object Ticket object dtype: object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Parch	int64
<pre>Pclass object Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv</pre>	PassengerId	object
Sex object SibSp int64 Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Pclass	object
Survived object Ticket object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	SibSp	int64
Ticket object dtype: object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Survived	object
dtype: object Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_train.csv	Ticket	object
Data types in test dataset after conversion: Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agent competition/titanic/cleaned_train.csv	dtype: object	
Age float64 Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agent competition/titanic/cleaned_train.csv	Data types in	test dataset after conversion.
Embarked object Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agent competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	Age	float64
Fare float64 Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agent competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	Embarked	object
Name object Parch int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agent competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	Fare	float64
Parson Int64 PassengerId object Pclass object Sex object SibSp int64 Ticket object dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	Name	object
Pclass object Sex object SibSp int64 Ticket object dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen	Parch	INT04
<pre>Sex object SibSp int64 Ticket object dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen</pre>	Pclass	object
<pre>SibSp int64 Ticket object dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen</pre>	Sex	object
<pre>Ticket object dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen</pre>	SibSp	int64
<pre>dtype: object Cleaned training data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agen</pre>	Ticket	object
<pre>cleaned training data saved to /mmt/d/rythonProjects/AutoKaggleMaster/multi_agen competition/titanic/cleaned_train.csv Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv</pre>	dtype: object	and data asynch to (mat/d/DythenDrof-st-("))t-V-st-V-st-V-st-V-
Cleaned test data saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ competition/titanic/cleaned_test.csv	competitio	ng data saved to /mnt/d/rythonrrojects/AutoKaggieMaster/multi_agen pn/titanic/cleaned train.csv
competition/titanic/cleaned_test.csv	Cleaned test d	ata saved to /mnt/d/PythonProjects/AutoKaggleMaster/multi agents/
	competitio	on/titanic/cleaned_test.csv

E.4 IN-DEPTH EXPLORATORY DATA ANALYSIS

1887

1888

1889 In this state, the AutoKaggle delves further into the Titanic dataset. 1) The process begins with an extended univariate analysis to explore the distribution of both numerical and categorical features.

1890	Key statistical summaries are generated for numerical features such as age and fare, while bar charts
1891	and frequency tables are created for categorical features like passenger class and gender. 2) A bi-
1892	variate analysis investigates the relationship between individual features and the survival outcome.
1893	Box plots and violin plots are used to analyze how numerical features vary between survivors and
1894	non-survivors, while count plots are generated for categorical features to visualize survival rates
1895	across different groups. 3) A correlation analysis is conducted to explore the relationships between
1896	numerical features, visualized through a correlation matrix and heatmap. This helps to identify any
1897	strong correlations between features and the target variable, survival. 4) A multivariate analysis is
1898	performed to explore interactions between key features such as passenger class, gender, and age,
1800	in relation to survival. Visualizations, such as stacked bar charts and facet grids, are used to high-
1000	light these complex interactions, providing deeper insights into the data.Below is an excerpt of the
1900	generated file:
1901	

Gı	reat! Based on the information provided, we will proceed with a detailed plan for **In-depth Exploratory Data Analysis (EDA)** phase. The aim is to perform a thereare analysis of each feature to uncover decome incident that can inform
	feature engineering and model building.
# 1	## Plan for In-Depth Exploratory Data Analysis
# 1	### Task 1: Detailed Univariate Analysis
	feature.
1	. **Age:**
	 - **Actions:** - Plot a histogram to observe the distribution.
	- Use a boxplot to identify any remaining outliers.
	- **Tools & Parameters:**
	<pre>- `seaborn.boxplot(data=cleaned_train, x='Age')`</pre>
	- **Expected Output:**
	- Boxplot indicating the presence of outliers.
	- **Constraints:**
	- Ensure plots are clear and readable.
2.	. **Fare:**
З.	. **Categorical Features (Pclass, Sex, Embarked):**
# #	### Task 2: Bivariate Analysis
* 7	*Objective:** Explore relationships between pairs of features, especially with the
	carget variable barvivea.
1.	. **Survived vs. Age:**
2.	. **Survived vs. Fare:**
3.	. $\star\star$ Survived vs. Categorical Features (Pclass, Sex, Embarked): $\star\star$
# #	### Task 3: Multivariate Analysis
* 1	Objective:** Understand interactions between multiple features and their combined effect on the target variable.
1.	. **Pairplot:**
2.	. **Correlation Heatmap:**
# #	### Task 4: Insights Summary
* '	vojective.** Summarize key finaings from the in-depth EDA to inform the next phas
1.	. **Actions:**
2	 . **Tools & Parameters:**
3.	. **Expected Output:**
4	**Constraints:**
B	y following this detailed plan, you will be able to gain a thorough understanding
	the dataset, which will be invaluable for the subsequent phases of feature
	engineering and model building.

In-depth Exploratory Data Analysis Code

1998

```
2000
             import pandas as pd
2001
             import matplotlib.pyplot as plt
2002
             import numpy as np
             # Load the cleaned training data
2004
             cleaned_train = pd.read_csv('cleaned_train.csv')
             # Task 1: Univariate analysis on the Age and Fare features
2006
            # Histogram and Boxplot for Age and Fare
for feature in ['Age', 'Fare']:
2007
                plt.figure(figsize=(10, 6))
2008
                 cleaned_train[feature].hist(bins=30)
                plt.title(f'{feature} Distribution')
2009
                plt.xlabel(feature)
2010
                plt.ylabel('Frequency')
                plt.savefig(f'deep_eda/images/{feature.lower()}_histogram.png')
2011
                plt.close()
2012
                plt.figure(figsize=(10, 6))
                 cleaned_train[feature].plot(kind='box')
2013
                plt.title(f'Boxplot of {feature}')
2014
                plt.xlabel(feature)
                plt.savefig(f'deep_eda/images/{feature.lower()}_boxplot.png')
2015
                plt.close()
2016
            print("Univariate analysis on Age and Fare features completed.")
2017
2018
             # Task 2: Univariate analysis on categorical features
             # Countplot for Pclass, Sex, Embarked
2019
             for feature in ['Pclass', 'Sex', 'Embarked']:
2020
                cleaned_train[feature].value_counts().plot(kind='bar', figsize=(10, 6), title=f'{
                 feature} Distribution')
2021
                plt.xlabel(feature)
2022
                plt.ylabel('Frequency')
                plt.savefig(f'deep_eda/images/{feature.lower()}_countplot.png')
2023
                plt.close()
2024
             # Task 3: Bivariate analysis of Survived vs. Age, Fare, categorical features
2025
             # Violin plot for Age and Fare by Survived
2026
             for feature in ['Age', 'Fare']:
                plt.figure(figsize=(10, 6))
2027
                cleaned_train.boxplot(column=feature, by='Survived')
2028
                plt.title(f'Survival by {feature}')
                plt.xlabel('Survived')
2029
                plt.ylabel(feature)
2030
                plt.savefig(f'deep_eda/images/survived_vs_{feature.lower()}.png')
                plt.close()
2031
2032
             # Countplot for categorical features by Survived
             for feature in ['Pclass', 'Sex', 'Embarked']:
2033
                pd.crosstab(cleaned_train[feature], cleaned_train['Survived']).plot(kind='bar',
2034
                 stacked=True, figsize=(10, 6))
                plt.title(f'Survival by {feature}')
2035
                plt.xlabel(feature)
2036
                plt.ylabel('Count')
                plt.savefig(f'deep_eda/images/survived_vs_{feature.lower()}.png')
2037
                plt.close()
2038
             # Task 4: Multivariate analysis using a correlation heatmap
             # Correlation heatmap
2040
             numeric_df = cleaned_train.select_dtypes(include=[np.number])
             plt.figure(figsize=(10, 8))
2041
            plt.matshow(numeric_df.corr(), cmap='coolwarm', fignum=1)
2042
             plt.title('Correlation Heatmap')
             plt.savefig('deep_eda/images/correlation_heatmap.png')
2043
            plt.close()
2044
             # Task 5: Summarize key insights from the EDA
2045
             summary = """
2046
             ....
....
2047
             # Save the summary to a text file
             with open('deep_eda/eda_summary.txt', 'w') as file:
                 file.write(summary)
2049
2050
```

In-dee	p EDA
Summar count mean	y statistics for Age: 891.000000 29.039282
std	12.072074 2.500000
25%	22.000000
50% 75%	28.000000
max	54.500000
Name:	Age, dtype: float64
Cumrier	al vata hu Dalaaa.
Pclas	ai face by refass: S
1 0	. 629630
3 0	.4/2826
Name:	Survived, dtype: float64
Surviv	al rate by Sex:
female	0.742038
male	0.188908
Surviv	al rate by Embarked:
Embar	ked
	.553571
S 0	.339009
Name:	Survived, dtype: float64
COLLET	Age SibSp Parch Fare Survived
Age	1.000000 -0.239601 -0.178959 0.144544 -0.060622
SibSp Parch	-0.239601 1.000000 0.414838 0.332021 $-0.035322-0.178959$ 0.414838 1.000000 0.292616 0.081629
Fare	0.144544 0.332021 0.292616 1.000000 0.317430
Surviv	ed -0.060622 -0.035322 0.081629 0.317430 1.000000

E.5 FEATURE ENGINEERING

In this phase, the AutoKaggle add several new features to enhance the predictive power of the dataset. 1) A FamilySize feature is created by combining the SibSp and Parch columns, representing the total number of family members aboard, including the passenger. This feature captures the familial context, which could influence survival likelihood. 2) An AgeGroup feature is derived by categorizing passengers into age groups, simplifying the continuous age variable into meaningful categories such as "Child" and "Senior." This transformation helps identify potential age-related survival patterns. 3) Categorical features like Sex, Embarked, and Pclass are then encoded into numerical form to ensure they can be used in the model. One-hot encoding is applied to Sex and Embarked, while label encoding is used for Pclass, respecting its ordinal nature. 4) The cabin data is processed by extracting the first letter of the Cabin feature to create a new Deck variable. This feature provides information about the passenger's location on the ship, which may correlate with survival outcomes. Missing cabin data is handled by assigning an 'Unknown' category, ensuring completeness of the feature.

2106	ſ
2107	
2108	
2109	
2110	
2111	
2112	
2113	
2114	
2115	
2116	
2117	
2118	
2119	
2120	
2121	I
2122	
2123	I
2124	
2125	
2126	
2127	
2128	
2129	
2130	
2131	
2132	
2133	
2134	
2135	
2136	
2137	
2138	
2139	
2140	
2141	I
2142	
2144	I
2144	
2146	
2147	I
2148	I
2149	I
2150	I
2151	I
2152	I
2153	I
2154	I
2155	I
2156	
2157	1

Great! Let curre speci:	's design a detailed plan for the Feature Engineering phase, focusing on the ent features and the available tools. We'll ensure that each task is clear, fic, and includes the necessary steps to achieve the desired outcome.
### Featur	e Engineering Plan
#### Task **Objectiv addit:	l: Create New Features e:** Enhance the dataset by creating new features that could capture ional information relevant to predicting survival.
1. **Famil - **Act - **Met	y Size:** ion:** Create a new feature `FamilySize` by combining `SibSp` and `Parch`. hod:** `FamilySize = SibSp + Parch + 1` (including the passenger themselves)
- **Imp can :	act:** Captures the total number of family members traveling together, which influence survival chances.
2. **IsAlo	ne:**
3. **Age B	ins:**
4. **Fare	per Person:**
#### Task **Objectiv learn	2: Encode Categorical Features e:** Convert categorical features into numerical format suitable for machine hing models.
1. **Sex:*	*
2. **Embar	ked.**
Objectiv are cu	e: Extract useful information from 'Ticket' and 'Cabin' features, which urrently in text format.
1. **Ticke	t:**
2. **Cabin	**
#### Task **Objectiv , impi	4: Scale Numerical Features e:** Standardize numerical features to ensure they are on a comparable scale roving model performance.
1. **Numer	ical Features:**
<pre>### Summar - **New Fe - **Encode - **Proces - **Scaled</pre>	y of Expected Output atures Added:** `FamilySize`, `IsAlone`, `AgeBins`, `FarePerPerson` d Features:** `Sex` (label encoded), `Embarked` (one-hot encoded) sed Features:** `TicketPrefix`, `CabinDeck` Features:** `Age`, `Fare`, `FarePerPerson`
### Constr - **Runtim	aints and Considerations e Efficiency:** Ensure feature creation and encoding steps are optimized for
- **Handli featur	ng Missing Values:** Address any missing values in the newly created res appropriately.
- **Avoid sets t	Data Leakage:** Perform encoding and scaling separately on train and test to prevent data leakage.
By followi and mo ensu	ng these tasks, you will transform the cleaned data into a more informative odel-ready format (`processed_train.csv` and `processed_test.csv`). This pla res that the most critical steps of feature engineering are covered,

Feature Engineering Code

2160

2161

```
2162
            import pandas as pd
2163
2164
            # Load the cleaned datasets
            train_df = pd.read_csv('cleaned_train.csv')
2165
            test_df = pd.read_csv('cleaned_test.csv')
2166
            # Create copies of the dataframes
2167
            train_df_copy = train_df.copy()
2168
            test_df_copy = test_df.copy()
2169
            # STEP 1: Create New Features
2170
            train_df_copy['FamilySize'] = train_df_copy['SibSp'] + train_df_copy['Parch'] + 1
            test_df_copy['FamilySize'] = test_df_copy['SibSp'] + test_df_copy['Parch'] + 1
2171
2172
            train_df_copy['IsAlone'] = (train_df_copy['FamilySize'] == 1).astype(int)
            test_df_copy['IsAlone'] = (test_df_copy['FamilySize'] == 1).astype(int)
2173
2174
            age_bins = [0, 12, 18, 35, 60, 120]
            age_labels = ['Child', 'Teen', 'Adult', 'Senior', 'Elder']
2175
            train_df_copy['AgeBins'] = pd.cut(train_df_copy['Age'], bins=age_bins, labels=
2176
                 age_labels, right=False)
            test_df_copy['AgeBins'] = pd.cut(test_df_copy['Age'], bins=age_bins, labels=age_labels
2177
                , right=False)
2178
            train_df_copy['FarePerPerson'] = train_df_copy['Fare'] / train_df_copy['FamilySize']
2179
            test_df_copy['FarePerPerson'] = test_df_copy['Fare'] / test_df_copy['FamilySize']
2180
            # Save the datasets with new features
2181
            train_df_copy.to_csv('processed_train.csv', index=False)
2182
            test_df_copy.to_csv('processed_test.csv', index=False)
2183
            # Load the processed datasets
2184
            train_df = pd.read_csv('processed_train.csv')
            test_df = pd.read_csv('processed_test.csv')
2185
2186
            # Define functions to extract ticket prefix and cabin deck
            def extract_ticket_prefix(ticket):
2187
                parts = ticket.split()
2188
                return parts[0] if not parts[0].isdigit() else 'None'
2189
            def extract_cabin_deck(cabin):
2190
                return cabin[0] if pd.notna(cabin) else 'Unknown'
2191
            # Extract TicketPrefix and CabinDeck
2192
            train_df['TicketPrefix'] = train_df['Ticket'].apply(extract_ticket_prefix)
            test_df['TicketPrefix'] = test_df['Ticket'].apply(extract_ticket_prefix)
2193
2194
            train_df['CabinDeck'] = train_df['Cabin'].apply(extract_cabin_deck) if 'Cabin' in
                 train_df.columns else 'Unknown'
2195
            test_df['CabinDeck'] = test_df['Cabin'].apply(extract_cabin_deck) if 'Cabin' in
2196
                 test_df.columns else 'Unknown'
2197
            # Save the datasets with extracted features
2198
            train_df.to_csv('processed_train.csv', index=False)
            test df.to csv('processed test.csv', index=False)
2199
2200
```

E.6 MODEL BUILDING, VALIDATION, AND PREDICTION

In this phase, we conduct a comprehensive analysis of the Titanic passenger dataset with the aim of predicting passengers' survival probabilities. Initially, the data undergo preprocessing that included filling missing values, deleting columns with high missingness, and handling outliers. Subsequent feature engineering efforts introduce new attributes such as family size, solitary travel, age groupings, and fare per person, and involved encoding for gender and embarkation points. Furthermore, a random forest model is employed, optimized via grid search, and evaluated using cross-validation. Predictions are then made on the test set, and a submission file is prepared.

2211

2212

2201 2202

2203

~~ ~ ~

2214 Model Building, Validation, and Prediction Plan 2215 2216 ### Detailed Plan 2217 2218 #### Task 1: Prepare Training Data **Objective:** Separate the target variable and remove non-numeric columns. 2219 2220 1. **Separate Target Column** - **Action:** Extract the 'Survived' column from 'processed_train.csv' as 'y'. 2221 - **Tool:** pandas 2222 - **Expected Output:** 'y' as a separate pandas Series containing the target variable. 2223 - **Code Example:** ```python 2224 import pandas as pd 2225 train_data = pd.read_csv('/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/ 2226 competition/titanic/processed_train.csv') y = train_data['Survived'] 2227 2228 2. **Remove Non-Numeric Columns** 2229 - **Action:** Identify and drop non-numeric columns from the training set. 2230 - **Features Involved:** 'Name', 'Ticket', 'TicketPrefix', 'CabinDeck' - **Tool:** pandas 2231 - **Expected Output:** 'X_train' as a DataFrame containing only numeric columns. 2232 - **Code Example:** '''python 2233 X_train = train_data.drop(columns=['Survived', 'Name', 'Ticket', 'TicketPrefix', 2234 'CabinDeck']) ... 2235 2236 #### Task 2: Prepare Test Data **Objective:** Ensure the test data has the same structure as the training data. 2237 2238 1. **Remove Non-Numeric Columns** 2239 2240 #### Task 3: Train and Validate the Model **Objective:** Train and validate a machine learning model using the prepared datasets 2241 2242 1. **Model Selection and Training** 2243 2244 2. **Hyperparameter Tuning** 2245 2246 #### Task 4: Make Predictions 2247 **Objective:** Use the trained model to make predictions on the test dataset. 2248 1. **Prediction** 2249 2250 #### Task 5: Prepare Submission File 2251 **Objective:** Prepare the submission file in the required format. 2252 1. **Submission File Preparation** 2253 2254 ### Final Notes 2255 - Ensure that all steps are executed in sequence to maintain data integrity. 2256 - Validate each step with print statements or logs to confirm successful execution. - Use cross-validation to ensure model robustness and avoid overfitting. 2257 2258 Feel free to execute the tasks based on this detailed plan. If you have any questions or need further adjustments, don't hesitate to ask. 2259 2260 2261 2262

Model Building, Validation, and Prediction Code

competition/titanic/processed_train.csv'

train_df = pd.read_csv(processed_train_path)

Identify non-numeric columns to be dropped

Drop the target variable and non-numeric columns

X_train = train_df.drop(columns=columns_to_drop)

processed_train_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/

columns_to_drop = ['Survived', 'Name', 'Ticket', 'TicketPrefix', 'CabinDeck', 'AgeBins

processed_test_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition

X_test = test_df.drop(columns=['Name', 'Ticket', 'TicketPrefix', 'CabinDeck', 'AgeBins

Load the processed training data

Separate the target variable

Load the processed test data

X_test = X_test[X_train.columns]

(1)

param grid = {

/titanic/processed_test.csv

test_df = pd.read_csv(processed_test_path)

Ensure the column order matches X_train

Define the model and parameter grid

grid_search.fit(X_train, y)

Get the best model

})

'n_estimators': [100, 200, 300], 'max_depth': [None, 10, 20, 30]

best_model = grid_search.best_estimator_
Validate the model using cross-validation

Make predictions on the test data
predictions = best_model.predict(X_test)

Convert predictions to a DataFrame
predictions_df = pd.DataFrame({

'Survived': predictions

Prepare the submission file

titanic/submission.csv'

'PassengerId': test_df['PassengerId'],

predictions_df.to_csv(submission_path, index=False)

print(f"Test data processed. X_test shape: {X_test.shape}")

from sklearn.model_selection import GridSearchCV, cross_val_score

grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy')

cv_scores = cross_val_score(best_model, X_train, y, cv=5, scoring='accuracy')

print(f"Predictions made. Sample predictions:\n{predictions_df.head()}")

submission_path = '/mnt/d/PythonProjects/AutoKaggleMaster/multi_agents/competition/

from sklearn.ensemble import RandomForestClassifier

model = RandomForestClassifier(random_state=42)

Use GridSearchCV for hyperparameter tuning

Drop non-numeric columns to match X_train

y = train_df['Survived']

import pandas as pd

```
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2287
2288
2289
2290
2291
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
```

2322 E.7 RESEARCH REPORT

After completing the entire data science pipeline, AutoKaggle reviews the interim reports from each phase and finally compiles a Research Report to document its key findings, major decisions, and their rationale at every phase. This helps users better understand the solutions provided by AutoKaggle.

2327 **Research Report** 2328 2329 # COMPETITION RESEARCH REPORT 2330 2331 ## 1. PRELIMINARY EDA 2332 The preliminary exploratory data analysis (EDA) involved examining the dataset's 2333 structure and identifying key features. The analysis revealed critical insights 2334 such as: **Missing Values**: Notably, the 'Age' and 'Cabin' features had significant missing 2335 values. **Outliers**: Outliers were identified in the 'Fare' and 'Age' features, with fares 2336 exceeding 200 and ages above 60. 2337 - **Statistical Summary**: The median age was around 30, while the median fare was approximately 14.45. 2338 2339 ## 2. DATA CLEANING 2340 2341 The data cleaning phase focused on addressing missing values and outliers. Key actions included: 2342 - **Imputation**: 2343 'Age': Missing values were filled with the mean. - 'Cabin': Filled with 'Missing' to indicate absence. 2344 - 'Embarked': Missing values filled with the mode. 2345 - **Outlier Treatment**: 'Fare' and 'Age' were clipped using the IQR method to reduce the influence of 2346 outliers. 2347 . . . 2348 ## 3. DEEP EDA 2349 In-depth EDA revealed valuable insights about the relationships between features and 2350 survival rates: 2351 - **Pclass**: Higher classes showed higher survival rates; average Pclass was 2.31. - **Sex**: Females had a significantly higher survival rate compared to males. 2352 - **Fare**: Higher fares correlated positively with survival, indicating a potential 2353 link to better accommodation and safety. . . . 2354 2355 ## 4. FEATURE ENGINEERING 2356 Feature engineering involved creating and transforming features to enhance model 2357 performance: - **New Features Created**: 2358 - **Title**: Extracted from 'Name'. 2359 - **FamilySize**: Summation of 'SibSp' and 'Parch'. - **IsAlone**: Binary feature indicating if a passenger traveled alone. 2360 - **FarePerPerson**: Calculated fare per individual in a family. 2361 . . . 2362 ## 5. MODEL BUILDING, VALIDATION, AND PREDICTION 2363 Multiple models were trained during this phase, including: 2364 - **Models**: XGBoost, SVM, Random Forest, Decision Tree, and Logistic Regression. 2365 - **Best Model**: Random Forest achieved the highest validation score of 0.8379. 2366 2367 ## 6. CONCLUSION 2368 The competition's approach involved a structured process of EDA, data cleaning, 2369 feature engineering, and model evaluation. Key insights included the strong influence of 'Sex', 'Pclass', and 'Fare' on survival rates. The most impactful 2370 decisions involved addressing missing values and outliers, which collectively 2371 improved data quality and model accuracy. Future recommendations include further feature engineering, hyperparameter tuning, and validation of feature importance 2372 to enhance model performance. 2373 2374