# DS-STAR: DATA SCIENCE AGENT VIA ITERATIVE PLANNING AND VERIFICATION

#### **Anonymous authors**

000

001

002 003 004

006

008

010 011

012

013

014

016

017

018

019

021

025

026

027 028

029

031

032

034

040

041

042

043 044

046

047

048

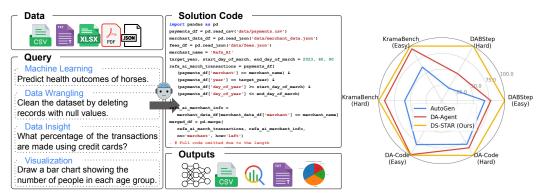
051

052

Paper under double-blind review

#### **ABSTRACT**

Data science, which transforms raw data into actionable insights, is critical for data-driven decision-making. However, these tasks are often complex, involving steps for exploring multiple data sources and synthesizing findings to deliver insightful answers. While large language models (LLMs) show significant promise in automating this process, they often struggle with heterogeneous data formats and generate sub-optimal analysis plans, as verifying plan sufficiency is inherently difficult without ground-truth labels for such open-ended tasks. To overcome these limitations, we introduce DS-STAR, a novel data science agent. Specifically, DS-STAR makes three key contributions: (1) a data file analysis module that automatically explores and extracts context from diverse data formats, including unstructured types; (2) a verification step where an LLM-based judge evaluates the sufficiency of the analysis plan at each stage; and (3) a sequential planning mechanism that starts with a simple, executable plan and iteratively refines it based on the DS-STAR's feedback until its sufficiency is verified. This iterative refinement allows DS-STAR to reliably navigate complex analyses involving diverse data sources. Our experiments show that DS-STAR achieves state-of-the-art performance across three challenging benchmarks: DABStep, KramaBench, and DA-Code. Moreover, DS-STAR particularly outperforms baselines on hard tasks that require processing multiple data files with heterogeneous formats.



(a) Data Science Tasks

(b) Experiment Summary

Figure 1: (a) Data science tasks require processing data files in various formats (e.g., csv, txt, xlsx, md) to answer user queries, such as data analysis for extracting useful insights or predictive tasks using machine learning. Data science agents are designed to accomplish this by writing code scripts (e.g., Python) and answering based on the output from their execution. In addition to the solution code, the output may include a trained model for prediction tasks, processed databases, text-formatted answers, visual charts, and more. (b) We report normalized accuracy (%) for both easy (answer can be found within a single file) and hard (requires processing multiple files) tasks in the DABStep, KramaBench, and DA-Code benchmarks. DS-STAR consistently and significantly outperforms competitive baselines, particularly in challenging hard tasks, showing DS-STAR's superiority in handling multiple data sources in heterogeneous formats.

# 1 Introduction

 The data science, which transforms raw data into actionable insights, is critical to solving real-world problems (De Bie et al., 2022; Hassan et al., 2023); for instance, businesses rely on insights extracted from data to make crucial directional decisions (Sun et al., 2018; Sarker, 2021). However, the data science workflow is often complex, requiring deep expertise across diverse fields such as computer science and statistics (Zhang et al., 2023; Hong et al., 2024). This workflow involves a series of time-consuming tasks (Zhang et al., 2024; Martin Iglesias, 2025), ranging from understanding distributed documents to intricate data processing and statistical analyses (see Figure 1).

To simplify this demanding workflow, recent research has explored using large language models (LLMs) as autonomous data science agents (Hong et al., 2024). These agents aim to translate natural language questions directly into functional code that generates corresponding answers (Ouyang et al., 2025; Wu et al., 2023; Yang et al., 2025; Yao et al., 2023). While representing significant advancements, current data science agents face several challenges that limit their practical utility and impact. A primary limitation is their predominant focus on well-structured data, such as relational databases composed of CSV files (Pourreza et al., 2024; Yu et al., 2018; Li et al., 2024). This narrow scope neglects the wealth information available in the heterogeneous data formats encountered in real-world scenarios, which can range from JSON to unstructured text and markdown files (Lai et al., 2025). Furthermore, as many data science tasks are framed as open-ended questions where ground-truth labels do not exist, verifying an agent's reasoning path is a non-trivial challenge. For example, most agents, like Data Interpreter (Hong et al., 2024), terminate their process upon successful code execution; however, executable code does not always guarantee a correct answer. Consequently, this lack of robust verification often leads an agent to adopt a sub-optimal plan for solving the given task.

To address these limitations, we introduce **DS-STAR**, a novel agent that tackles data science problems by generating solution code through a robust, iterative process of planning and verification. The proposed DS-STAR framework consists of two key stages. First, to ensure adaptability across diverse data types, DS-STAR automatically analyzes all files within a given directory, generating a textual summary of their structure and content. This summary then serves as a crucial context for the DS-STAR's approach to the given task. Second, DS-STAR enters a core loop of planning, implementation, and verification. It begins by formulating a high-level plan, implementing it as a code script, and then verifying to solve the problem. Notably, for verification, we propose using an LLM-based judge (Gu et al., 2024; Zheng et al., 2023) that is prompted to explicitly evaluate whether the current plan is sufficient to solve the problem. If the verification fails (*i.e.*, the judge determines the plan insufficient), DS-STAR refines its plan by adding or modifying steps and repeats the process. Crucially, rather than creating a complete plan at once, our agent operates sequentially, reviewing the results from each intermediate step before building upon them. <sup>1</sup> This iterative loop continues until a satisfactory plan is verified, at which point the corresponding code is returned as the final solution.

We validate the effectiveness of our proposed DS-STAR through comprehensive evaluations on a suite of established data science benchmarks. These benchmarks, including DABStep benchmark (Martin Iglesias, 2025), KramaBench (Lai et al., 2025), and DA-Code (Huang et al., 2024) are specifically designed to assess performance on complex data analysis tasks (*e.g.*, data wrangling, machine learning, visualization, exploratory data analysis, data manipulation, etc) involving multiple data sources and formats (see Section 4). The experimental results demonstrate that DS-STAR significantly outperforms existing state-of-the-art methods across all evaluated scenarios. Specifically, compared to the best alternative, DS-STAR improves accuracy from 41.0% to 45.2% on the DABStep benchmark, from 39.8% to 44.7% on KramaBench, and from 37.0% to 38.5% on DA-Code.

The contributions of this paper can be summarized as follows:

- We introduce DS-STAR, a novel data science agent which operates on data with various formats and on various tasks like data analysis, machine learning, visualization, and data wrangling, etc.
- We propose an LLM-based verification module that judges the sufficiency of a solution plan.
- We develop an iterative refinement strategy, enabling DS-STAR to sequentially build its solution.
- We provide a comprehensive evaluation on challenging benchmarks, showing DS-STAR's efficacy.

<sup>&</sup>lt;sup>1</sup>This mirrors the interactive process of an expert using a Jupyter notebook to observe interim results.

# 2 RELATED WORK

LLM agents. Recent advancements in LLMs have spurred significant research into autonomous agents designed to tackle complex and long-horizon tasks. A key strategy employed by these multi-agent systems is the autonomous decomposition of a main objective into a series of smaller, manageable sub-tasks, which can be delegated to sub-agents. General-purpose agents, such as ReAct (Yao et al., 2023), HuggingGPT (Shen et al., 2023), and OpenHands (Wang et al., 2024), utilize external tools to reason, plan, and act across a wide range of problems. Building on these foundational capabilities, research has increasingly focused on specialized domains such as Voyager (Wang et al., 2023) for navigating the Minecraft environment and AlphaCode (Li et al., 2022) for performing sophisticated code generation. Furthermore, DS-Agent (Guo et al., 2024), AIDE (Jiang et al., 2025), and MLE-STAR (Nam et al., 2025) are tailored specifically for machine learning engineering. In a similar vein, our work, DS-STAR, is a specialized LLM agent for data science tasks.

Data science agents. Recent research has focused on developing data science agents that leverage the advanced coding and reasoning capabilities of LLMs (Jiang et al., 2025; Jimenez et al., 2024). Initial efforts utilized general-purpose frameworks like ReAct (Yao et al., 2023) and AutoGen (Wu et al., 2023). Following these pioneering approaches, agents like DA-Agent (Huang et al., 2024), which are specialized in autonomously generating code-based solutions for tasks such as data analysis and visualization (Hu et al., 2024; Jing et al., 2025), has been developed. A notable example is Data Interpreter (Hong et al., 2024), which employs a graph-based method to decompose a primary task into a series of manageable sub-tasks and the resulting task graph is then progressively refined based on the successful execution of these sub-tasks. However, a key limitation is that relying solely on successful execution as feedback often results in sub-optimal plans, since such feedback cannot confirm the plan's correctness. To address this, DS-STAR introduces a novel verification process that employs an LLM as a judge (Gu et al., 2024) to assess the quality of the generated solutions.

**Text-to-SQL.** The task of Text-to-SQL involves translating natural language questions into executable SQL queries (Li & Jagadish, 2014; Li et al., 2023). Early approaches relied on sequence-to-sequence architectures, which jointly encoded the user's query and the database schema (Cao et al., 2021; Choi et al., 2021). However, after the rise of LLMs, initial research focused on prompt engineering (Pourreza & Rafiei, 2023). This was soon followed by the development of more sophisticated, multi-step pipelines incorporating techniques such as schema linking, self-correction, and self-consistency to enhance generation accuracy (Talaei et al., 2024; Pourreza et al., 2024; Maamari et al., 2024). While these pipelines function as specialized agents for data analysis, their fundamental reliance on generating SQL restricts their use to structured, relational databases. In contrast, DS-STAR addresses this limitation by adopting a Python language, thereby framing our approach as a Text-to-Python task. To further facilitate this, we introduce a novel data file summarization mechanism which enables our method to demonstrate significantly broader applicability than traditional Text-to-SQL systems, as it can operate on a wide array of data formats, including JSON, Markdown, and unstructured text.

# 3 DS-STAR

In this section, we introduce DS-STAR, a framework for data science agents that leverages the coding and reasoning capabilities of LLMs to tackle data science tasks. In a nutshell, our approach first analyzes heterogeneous input data files to extract essential information (Section 3.1). Subsequently, DS-STAR generates a code solution through an iterative process (Section 3.2). The prompts and algorithms are detailed in Appendix G and A, respectively.

**Problem setup.** Our goal is to automatically generate a code solution s (i.e., a Python script) that answers query q using a given data files  $\mathcal{D}$ . We formulate this as a search problem over the space of all possible scripts, denoted by  $\mathcal{S}$ . The quality of any script  $s \in \mathcal{S}$  is evaluated by a scoring function h(s), which measures the correctness of its output,  $s(\mathcal{D})$ , against a ground-truth answer (e.g., accuracy) where  $s(\mathcal{D})$  represents the output generated by executing s, which loads the data  $\mathcal{D}$  to produce an answer. The objective is to identify the  $s^*$  that is:  $s^* = \arg\max_{s \in \mathcal{S}} h(s(\mathcal{D}))$ .

In this paper, we propose a multi-agent framework,  $\mathcal{A}$ , designed to process the query q and data files  $\mathcal{D}$ , which may be numerous and in heterogeneous formats. This framework is composed of n specialized LLM agents,  $\{\mathcal{A}_i\}_{i=1}^n$ , where each agent possesses a distinct functionality, as detailed in the subsequent sections (see Figure 2 for the problem setup and our agent's overview).

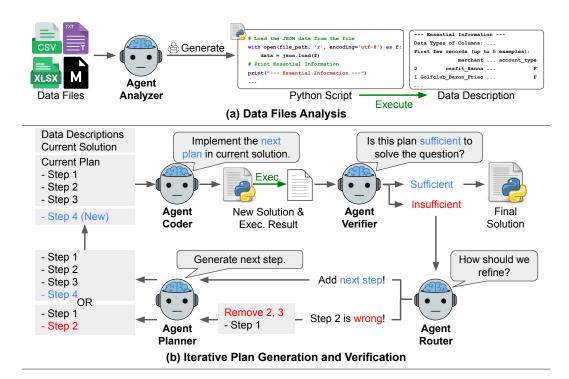


Figure 2: **Overview of DS-STAR.** (a) DS-STAR generates a Python script that analyzes heterogeneous input data files. (b) Starting with a simple single step, DS-STAR implements the plan and executes it to obtain intermediate results. Next, using a verifier agent, DS-STAR determines whether the current plan is sufficient to solve the user's query. If it is not sufficient, DS-STAR's router agent decides whether to add the next step or remove the incorrect step from the plan. Finally, DS-STAR adds the next step and implements it again. This iterative process of planning, implementation, execution, and verification is repeated until the verifier agent determines that the plan is sufficient to answer the user's query or until the maximum number of iterations is reached.

#### 3.1 Analyzing data files

To effectively interact with the given data files, DS-STAR first requires a comprehensive understanding of its contents and structure. We achieve this by generating a concise, analytical description for each data file. This initial analysis phase is critical as it informs all subsequent actions taken by the agent.

Conventional frameworks for data science agents often rely on displaying a few sample rows from structured files, such as CSVs. However, this approach is fundamentally limited to structured data and fails for unstructured formats where the concept of a 'row' is ill-defined. To overcome this limitation, we introduce a more general mechanism as follows. For each data file  $\mathcal{D}_i \in \mathcal{D}$  we employ an analyzer agent,  $\mathcal{A}_{\mathtt{analyzer}}$ , to generate a Python script,  $s_{\mathtt{desc}}^i$ . This script is designed to correctly parse and load the data file  $\mathcal{D}_i$  and then extract its essential properties. For structured data, this might include column names and data types; for unstructured data, it could be file metadata, text summaries. The resulting analytical description,  $d_i$ , is captured directly from the script's execution output (see Appendix D). This process, which can be parallelized, is denoted as:  $d_i = \mathsf{exec}(s_{\mathtt{desc}}^i)$ ,  $s_{\mathtt{desc}}^i = \mathcal{A}_{\mathtt{analyzer}}(\mathcal{D}_i)$ .

#### 3.2 ITERATIVE PLAN GENERATION AND VERIFICATION

**Initialization.** After generating analytic descriptions of the N data files, DS-STAR begins the solution generation process. First, a planner agent,  $\mathcal{A}_{\texttt{planner}}$ , generates an initial high-level executable step  $p_0$  (e.g., loading a data file) using the query q and obtained data descriptions:  $p_0 = \mathcal{A}_{\texttt{planner}}(q, \{d_i\}_{i=1}^N)$ .

This single executable step is then implemented as a code script  $s_0$  by a coder agent  $\mathcal{A}_{\mathtt{coder}}$  and the initial execution result  $r_0$  is then obtained by executing the script  $s_0$  as follows:

$$s_0 = \mathcal{A}_{coder}(p_0, \{d_i\}_{i=1}^N), \ r_0 = exec(s_0).$$
 (1)

**Plan verification.** The main challenge in data science tasks is guiding the refinement of a solution, as determining its correctness is often non-trivial, since there are no ground-truth label. To address this, we use an LLM as a judge to assess whether the current plan is sufficient for the user's query.

Our approach introduces a verifier agent  $\mathcal{A}_{\text{verifier}}$ . At any given round k in the problem-solving process, this agent evaluates the state of the solution, *i.e.*, whether the plan is sufficient to solve the problem. The evaluation is based on the cumulative plan  $p = \{p_0, \dots, p_k\}$ , the user's query q, the current solution code  $s_k$ , which is an implementation of the cumulative plan, and its execution result  $r_k$ . The operation of  $\mathcal{A}_{\text{verifier}}$  is denoted as follows:  $v = \mathcal{A}_{\text{verifier}}(p, q, s_k, r_k)$ .

Here, the output v is a binary variable: sufficient or insufficient. Note that our method does not just compare the plan to the query. By conditioning the judgement on  $s_k$  and its execution output  $r_k$ ,  $\mathcal{A}_{\text{verifier}}$  can provide more grounded feedback since it assesses whether  $s_k$  is well-implemented following the plan, and whether the  $r_k$  contains the information needed to fully address the query.

**Plan refinement.** If the verifier agent  $\mathcal{A}_{\text{verifier}}$  determines that the current plan is insufficient to solve the user's query q, DS-STAR must decide how to proceed. Such insufficiency could arise because the plan is merely incomplete and requires additional steps, or because it contains erroneous steps that invalidate the approach. To resolve this, DS-STAR employs a router agent,  $\mathcal{A}_{\text{router}}$ , which decides whether to append a new step or to correct an existing one. The router's decision w is generated as follows, where  $p = \{p_0, \cdots, p_k\}$  is the current cumulative plan:  $w = \mathcal{A}_{\text{router}}(p, q, r_k, \{d_i\}_{i=1}^N)$ .

The output w is either the token Add Step or an index  $l \in \{1, \cdots, k\}$ . If w is Add Step,  $\mathcal{A}_{\mathtt{router}}$  has determined the plan is correct but incomplete. In this case, we retain the plan p and proceed to generate the next step. On the other hand, if w = l,  $\mathcal{A}_{\mathtt{router}}$  has identified  $p_l$  as erroneous. In this case, we backtrack by truncating the plan to  $p \leftarrow \{p_0, \cdots, p_{l-1}\}$ . Here, we deliberately choose to truncate and regenerate through the LLM's random sampling, rather than directly correcting  $p_l$ , since our empirical finding has shown that revising a specific incorrect step often leads to an overly complex replacement, therefore frequently flagged again by  $\mathcal{A}_{\mathtt{router}}$  in a next iteration. Following the decision from the  $\mathcal{A}_{\mathtt{router}}$ , our agent proceeds with an updated plan  $p = \{p_0, \cdots, p_{k'}\}$ , where k' = k or k' = l - 1. Then DS-STAR generates a subsequent step:  $p_{k'+1} = \mathcal{A}_{\mathtt{planner}}(p, q, r_k, \{d_i\}_{i=1}^N)$ .

Notably, the planner agent  $\mathcal{A}_{\texttt{planner}}$  is conditioned on the last execution result,  $r_k$ , enabling it to generate a step that attempts to resolve the previously identified insufficiency. Once the new step  $p_{k'+1}$  is defined, the plan is updated to:  $p \leftarrow \{p_0, \cdots, p_{k'}, p_{k'+1}\}$ .

**Plan implement and execution.** Finally, DS-STAR enters an execution and verification cycle. First,  $\mathcal{A}_{\texttt{coder}}$  implements p into code s. The execution of this code yields a new observation r = exec(s). With this r,  $\mathcal{A}_{\texttt{verifier}}$  is invoked again to assess if the newly augmented plan is now sufficient. This entire iterative procedure—routing, planning, coding, executing, and verifying—is repeated until  $\mathcal{A}_{\texttt{verifier}}$  returns a sufficient or a predefined maximum number of iterations is reached.

#### 3.3 Additional modules for robust data science agents

**Debugging agent.** When a Python script s fails during execution, it generates an error traceback  $\mathcal{T}_{\text{bug}}$ . To automatically debug the script, DS-STAR employs a debugging agent,  $\mathcal{T}_{\text{debugger}}$ . First, when generating  $\{d_i\}_{i=1}^N$  using  $s_{\text{desc}}$  obtained from  $\mathcal{A}_{\text{analyzer}}$  (see Section 3.1),  $\mathcal{A}_{\text{debugger}}$  iteratively update the script using only the traceback:  $s_{\text{desc}} \leftarrow \mathcal{A}_{\text{debugger}}(s_{\text{desc}}, \mathcal{T}_{\text{bug}})$ .

Secondly, once DS-STAR obtains  $\{d_i\}_{i=1}^N$ ,  $\mathcal{A}_{\text{debugger}}$  utilizes such information when generating a solution Python script s (see Section 3.2). Our key insight is that tracebacks alone are often insufficient for resolving errors in data-centric scripts, while  $\{d_i\}_{i=1}^N$  might include critical metadata such as column headers in a CSV file, sheet names in an Excel workbook, or database schema information. Therefore,  $\mathcal{A}_{\text{debugger}}$  generates a corrected script, s, by conditioning on the original script s, the error traceback  $\mathcal{T}_{\text{bug}}$ , and this rich data context  $\{d_i\}_{i=1}^N$ :  $s \leftarrow \mathcal{A}_{\text{debugger}}(s, \mathcal{T}_{\text{bug}}, \{d_i\}_{i=1}^N)$ .

**Retriever.** A potential scalability challenge arises when the number of data files N is large (i.e., N > 100.). In such cases, the complete set of descriptions  $\{d_i\}_{i=1}^N$  cannot be prompted within the predefined context length of LLMs. To address this limitation, we employ a retrieval mechanism that leverages a pre-trained embedding model (Nie et al., 2024). Specifically, we identify the top-K most relevant data files, which will be provided as context to the LLM, by computing the cosine similarity between the embedding of the user's query q and the embedding of each description  $d_i$ .

Table 1: **Main results from DABStep.** All results are taken from the DABStep leaderboard (Martin Iglesias, 2025), except for the model marked with †. The highest scores are shown in **bold**.

Framework	Model	Easy-level Accuracy (%)	Hard-level Accuracy (%)
Model-only	Gemini-2.5-Pro	66.67	12.70
	o4-mini	76.39	14.55
ReAct (Yao et al., 2023)	Claude-4-Sonnet	81.94	19.84
	Gemini-2.5-Pro <sup>†</sup>	69.44	10.05
AutoGen (Wu et al., 2023)	Gemini-2.5-Pro <sup>†</sup>	59.72	10.32
Data Interpreter (Hong et al., 2024)	Gemini-2.5-Pro <sup>†</sup>	72.22	3.44
DA-Agent (Huang et al., 2024)	Gemini-2.5-Pro <sup>†</sup> DeepSeek-V3 Claude-3.5-Sonnet	68.06	22.49
Open Data Scientist		84.72	16.40
Mphasis-I2I-Agents		80.56	28.04
Amity DA Agent DS-STAR (Ours)	Gemini-2.5-Pro	80.56	41.01
	Gemini-2.5-Pro <sup>†</sup>	<b>87.50</b>	<b>45.24</b>

# 4 EXPERIMENTS

In this section, we conduct a comprehensive empirical evaluation of DS-STAR on three challenging data science benchmarks: DABStep (Martin Iglesias, 2025), KramaBench (Lai et al., 2025), and DA-Code (Huang et al., 2024). Our results show that DS-STAR significantly outperforms competitive baselines, including those built upon various LLMs (Section 4.1). To understand the source of gains, we perform a detailed ablation study, which confirms that our proposed analyzer and router agents are critical components that provides a substantial performance improvement (Section 4.2).

**Common setup.** Unless otherwise specified, all experiments are conducted using Gemini-2.5-Pro as the base LLM. Our agent, DS-STAR, operates for a maximum of 20 rounds per task. To ensure properly formatted output, we employ a finalyzer agent,  $\mathcal{A}_{\texttt{finalyzer}}$ , which takes formatting guidelines (e.g., rounding to two decimal places) and generates the final solution code (see Appendix G).

#### 4.1 Main results

**DABStep.** We evaluate DS-STAR on the DABStep benchmark (Martin Iglesias, 2025), which is designed to mirror real-world data analysis challenges by requiring the processing of seven diverse data files, including formats like JSON, Markdown, and CSV. The benchmark features two difficulty levels; its hard-level tasks are particularly demanding, necessitating the analysis of multiple data files and the application of domain-specific knowledge. As shown in Table 1, existing single-LLM approaches struggle to surpass 20% accuracy, which underscores the significant room for improvement. Additionally, its authors note that a human-reference solution for a hard-level task spans 220 lines of code and is broken into four sequential steps. Detailed statistics and the full execution logs of DS-STAR are provided in Appendix B and F, respectively.

The results, presented in Table 1, demonstrate that DS-STAR significantly outperforms all baselines. For instance, integrating DS-STAR with Gemini-2.5-Pro boosts the hard-level accuracy from 12.70% to 45.24%, an absolute improvement of over 32 percentage points. Notably, the DS-STAR using Gemini-2.5-Pro substantially surpasses other commercial agents like Open Data Scientist, Mphasis-I2I-Agents, and Amity DA Agent across both easy and hard difficulty levels. Crucially, DS-STAR also outperforms other multi-agent systems built upon the same Gemini-2.5-Pro, highlighting that the performance gains stem directly from the proposed effective sub-agent orchestration.

**KramaBench.** We additionally validate DS-STAR's capability on Kramabench (Lai et al., 2025), a benchmark that requires data discovery, *i.e.*, selecting the correct data files from a vast data lake to address a user's query. For example, the Astronomy domain includes over 1,500 files, while only a small subset is relevant to any single query (see Appendix B for the detailed statistics). To navigate this, DS-STAR integrates a retrieval module that identifies the top 100 candidate files <sup>2</sup> based on embedding similarity between the user's query and data descriptions, computed with Gemini-Embedding-001 (Lee et al., 2025) (details can be found in Section 3.3).

<sup>&</sup>lt;sup>2</sup>If the total data is less than 100, we fully utilized all the data for the task.

Table 2: **Main results from KramaBench.** All results are taken from the original paper (Lai et al., 2025), except the results with Gemini-2.5-Pro. The highest scores are shown in **bold**.

Framework	Model	Domains						
ггашемогк	Model	Archaeology	Astronomy	Biomedical	Environment	Legal	Wildfire	Total
Original experimental setting for which the relevant data must be retrieved.								
	03	25.00	1.73	3.50	1.35	3.35	24.87	9.64
Model-only	GPT-40	0.00	1.41	1.98	0.45	1.46	1.45	1.62
	Claude-3.5-Sonnet	16.67	1.62	2.87	1.17	7.33	13.63	7.45
	03	25.00	3.53	8.95	19.60	13.89	50.73	22.08
DS-GURU	GPT-40	16.67	2.76	8.97	2.60	2.80	17.18	8.28
	Claude-3.5-Sonnet	16.67	1.52	1.96	11.21	7.01	39.16	14.35
ReAct	Gemini-2.5-Pro	16.67	4.77	3.69	26.17	41.95	51.40	30.31
AutoGen	Gemini-2.5-Pro	16.67	4.39	7.25	19.38	26.38	41.76	22.83
Data Interpreter	Gemini-2.5-Pro	41.67	12.72	28.05	9.87	30.04	59.67	31.32
DA-Agent	Gemini-2.5-Pro	41.67	15.52	12.59	42.64	39.73	61.61	39.79
DS-STAR (Ours)	Gemini-2.5-Pro	25.00	12.09	43.74	46.75	49.64	65.94	44.69
	Oracle ex	xperimental setti	ing with releva	nt data alread	y provided.			
ReAct	Gemini-2.5-Pro	25.00	6.36	3.69	30.70	46.93	51.69	33.82
AutoGen	Gemini-2.5-Pro	25.00	5.29	26.03	23.46	39.28	49.41	31.77
Data Interpreger	Gemini-2.5-Pro	41.67	13.29	28.39	13.09	32.97	63.13	33.57
DA-Agent	Gemini-2.5-Pro	41.67	15.77	44.26	44.55	56.42	65.92	48.61
DS-STAR (Ours)	Gemini-2.5-Pro	25.00	19.08	55.24	58.80	59.66	70.18	52.55

Table 3: **Main results from DA-Code.** All results are taken from the original paper (Huang et al., 2024), except the results with Gemini-2.5-Pro. The highest scores are shown in **bold**.

E	M - 1-1	Score						
Framework	Model	Data Wrangling	ML	EDA	Easy	Medium	Hard	Total
	Mixtral-8x22B	14.8	31.6	10.2	17.6	16.8	8.6	15.4
	DeepSeek-Coder-V2.5	25.1	34.1	14.7	32.8	18.7	14.1	20.7
DA A	Qwen-2.5-72B	24.9	41.8	15.4	31.9	19.4	22.3	22.6
DA-Agent	Claude-3-Opus	29.3	46.8	20.7	44.7	23.8	19.0	27.6
	GPT-4	30.4	48.4	24.6	45.4	27.8	23.4	30.5
	Gemini-2.5-Pro	34.8	57.2	31.1	50.0	34.2	32.0	37.0
ReAct	Gemini-2.5-Pro	14.7	31.2	22.2	32.0	17.9	25.9	22.5
AutoGen	Gemini-2.5-Pro	25.6	51.6	25.6	38.7	28.5	29.6	30.8
DS-STAR (Ours)	Gemini-2.5-Pro	30.4	57.3	34.8	48.9	35.2	37.1	38.5

As detailed in Table 2, DS-STAR demonstrates a substantial performance advantage in this challenging setting. Using Gemini-2.5-Pro, DS-STAR achieves an accuracy of 44.69%, significantly outperforming the 39.79% score of the state-of-the-art DA-Agent. Moreover, to isolate the impact of data retrieval, we also evaluated performance in an oracle setting, assuming all relevant data for a task is already provided. Under these ideal conditions, DS-STAR's accuracy increases by 8 percentage points. In contrast, other agentic frameworks like ReAct show only a marginal performance gain in the same setting. This large gap highlights that while our current retrieval method is effective, advanced data discovery is a promising direction for unlocking the full potential of DS-STAR.

**DA-Code.** To evaluate the generalization capabilities of DS-STAR, we additionally benchmark its performance on the DA-Code dataset (Huang et al., 2024), which encompasses a diverse array of data science tasks. Specifically, DA-Code is structured into three main categories: data wrangling, machine learning (ML), and exploratory data analysis (EDA), with the EDA category further subdivided into data manipulation, data insights, visualization, and statistical analysis. Here, we found that the Data-Interpreter (Hong et al., 2024) struggles to generate answers with the right format required by DA-Code tasks (*e.g.*, saving in csv format), therefore excluded it from the baseline.

As shown in Table 3, DS-STAR outperforms the strongest baseline, *i.e.*, DA-Agent, demonstrating its robust applicability across various data science domains. The superiority of our framework is particularly highlighted on more complex problems. Specifically, on the hard-difficulty tasks, DS-STAR achieves 37.1% accuracy, representing a significant improvement over the 32.0% accuracy of the DA-Agent baseline, even when both frameworks utilize the same Gemini-2.5-Pro. This showcases the effectiveness of our approach, especially for challenging, real-world data science scenarios. We provide qualitative results for each task type in Appendix E.

# 378 Ques

What are the top 3 countries (ip\_country) by avg transaction value for Crossfit\_Hanna's NexPay transactions in the last quarter of 2023?

```
DS-STAR (Ours)
import pandas as pd
payments_file = 'data/payments.csv
payments_df = pd.read_csv(payments_file)
# Filter the DataFrame for 'C
        and Q4 transa
filtered_transactions = payments_df[
   (payments_df['merchant'] == 'Crossfit_Hanna') &
   (payments df['card scheme'] == 'NexPay') &
   (payments df['day of year'] >= 274)
top_3_countries_by_avg_amount = (
   filtered transactions.groupby('ip country')['eur amount']
   .mean()
   .sort_values(ascending=False)
   .head(3)
answer = ','.join(top_3_countries_by_avg_amount.index)
print(answer)
```

Figure 3: **Qualitative result.** We observed that while ReAct (Yao et al., 2023) failed to filter out the 'NextPay' value, leading to an incorrect answer, DS-STAR successfully filtered out due to the use of data file analysis agent. See Appendix F for the full execution logs of DS-STAR.

Table 4: **Ablation study on each component.** Effectiveness of each agent and compatibility with other LLMs on the DABStep benchmark (Martin Iglesias, 2025).

Framework	Model	Analyzer $\mathcal{A}_{ ext{analyzer}}$	$\begin{array}{c} \textbf{Router} \\ \mathcal{A}_{\mathtt{router}} \end{array}$	Easy-level Accuracy (%)	Hard-level Accuracy (%)
DS-STAR (Variant 1)	Gemini-2.5-Pro	X	✓	75.00	26.98
DS-STAR (Variant 2)	Gemini-2.5-Pro	✓	×	79.17	39.95
DS-STAR (Ours)	Gemini-2.5-Pro	✓	✓	87.50	45.24
DS-STAR (Ours)	GPT-5	✓	✓	88.89	43.12

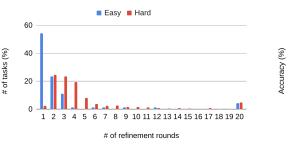
# 4.2 Ablation studies

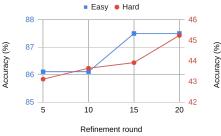
In this section, we conduct ablation studies to verify the effectiveness of individual components. Unless otherwise specified, all experiments use Gemini-2.5-Pro and are performed on the DABStep.

Effectiveness of the data file analysis agent. The data file analysis agent,  $\mathcal{A}_{\mathtt{analyzer}}$ , which provides contextual descriptions of data files, is crucial for achieving high performance. When we remove the descriptions, DS-STAR's accuracy on hard-level tasks in the DABStep benchmark drops significantly to 26.98%. While this result is still substantially better than the 12.70% accuracy of a non-agentic framework, it underscores the importance of the component. We hypothesize that the rich context of the given data is vital for enabling DS-STAR to effectively plan and implement its approach.

To qualitatively demonstrate the effectiveness of our data file analysis agent, we present a comparative case study in Figure 3. In this example, the competing baseline, ReAct (Yao et al., 2023), tries to filter data based on a 'NextPay' value but fails to execute the request. The model incorrectly determines that no such column value exists in the 'payment.csv' file due to ReAct's limited understanding of the provided data context. On the other hand, DS-STAR succeeds by first employing its data file analysis to correctly understand the entire data structure. This allows it to identify the appropriate column value and accurately filter the data as instructed. Further qualitative examples, including full execution logs, showcasing the capabilities of DS-STAR can be found in Appendix F.

Effectiveness of the router agent.  $\mathcal{A}_{router}$  determines whether to append a new step or correct an existing one in the plan. Here, we verify its importance against an alternative approach (Variant 2 in Table 4) where  $\mathcal{A}_{router}$  is removed. In this variant, DS-STAR only uses  $\mathcal{A}_{planner}$  to simply add new steps sequentially until the plan is sufficient or a maximum number of iterations is reached. As shown in Table 4, this alternative performs worse on both easy and hard tasks. This is because building upon an erroneous step leads to error accumulation, causing subsequent steps to also be incorrect, *i.e.*, correcting errors in the plan is more effective than accumulating potentially flawed steps.





(a) Distribution of number of required rounds

(b) Sensitivity analysis on number of refinements

Figure 4: **Sensitivity analysis on the number of refinement steps using DABStep.** (a) Difficult tasks require more iterations to generate sufficient plans. Specifically, hard-level tasks require an average 5.6 iterations, while easy-level tasks require 3.0 iterations. Also, more than 50% of easy-level tasks are done only with a single round. (b) Performing more iterations allows the agent to generate sufficient plans, resulting in better performance for both easy and hard level tasks.

**Compatibility with other LLMs.** To evaluate the generalizability of DS-STAR across LLM types, we conduct additional experiments using GPT-5. As shown in Table 4, DS-STAR with GPT-5 also achieves promising results on the DABStep benchmark. Notably, DS-STAR with GPT-5 demonstrates stronger performance on easy-level tasks, whereas the Gemini-2.5-Pro excels on hard-level tasks.

#### 4.3 ANALYSIS

In this section, we analyze the impact of the number of refinement rounds. Specifically, we measure the number of iterations required to generate a sufficient plan and conduct a sensitivity analysis on the maximum number of iterations. All experiments were done on the DABStep using Gemini-2.5-Pro. In addition, we provide cost analysis of DS-STAR in the Appendix C.

**Required number of refinement steps.** As shown in Figure 4(a), generating a successful plan for hard tasks requires a greater number of refinement rounds. Specifically, on the DABStep, hard tasks required an average of 5.6 rounds, in contrast to the 3.0 rounds needed for easy tasks. This disparity is further underscored by the fact that while over 50% of easy-level tasks were solved by the initial plan  $p_0$  alone, nearly all, *i.e.*, 98% of the hard tasks necessitated at least one refinement iteration.

Sensitivity analysis on the maximum number of refinement steps. In our experiments, DS-STAR uses a default maximum of 20 refinement iterations. To analyze the sensitivity of this hyper-parameter, we conduct an experiment by lowering the limits to 5, 10, and 15. If the process reaches this maximum without finding a sufficient plan, DS-STAR generates a final solution using the intermediate plan it has developed. Figure 4(b) shows a positive correlation between the maximum number of refinements and task accuracy for both easy and hard-level tasks. A higher iteration limit increases the probability that DS-STAR will generate a sufficient plan. This seems to be more critical for hard-level tasks, where accuracy consistently improves as the maximum number of rounds increases, verifying the importance of sufficient number of refinement steps for complex problems.

### 5 CONCLUSION

We introduce DS-STAR, a novel agent designed to autonomously solve data science problems. Our approach features two key components: (1) the automatic analysis of files to handle heterogeneous data formats, and (2) the generation of a sequential plan that is iteratively refined with a novel LLM-based verification mechanism. We show the effectiveness of DS-STAR on the DABStep, Kramabench, and DA-Code, where it establishes a new state-of-the-art by outperforming prior methods.

**Limitation and future works.** Our current work focuses on a fully automated framework for DS-STAR. A compelling avenue for future research is to extend this framework to a human-in-the-loop setting. Investigating how to synergistically combine the automated capabilities of DS-STAR with the intuition and domain knowledge of a human expert presents a promising direction for significantly boosting performance and enhancing the system's practical utility.

#### REFERENCES

- Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu, and Kai Yu. Lgesql: line graph enhanced text-to-sql model with mixed local and non-local relations. *arXiv preprint arXiv:2106.01093*, 2021.
- DongHyun Choi, Myeong Cheol Shin, EungGyun Kim, and Dong Ryeol Shin. Ryansql: Recursively applying sketch-based slot fillings for complex text-to-sql in cross-domain databases. *Computational Linguistics*, 2021.
  - Tijl De Bie, Luc De Raedt, José Hernández-Orallo, Holger H Hoos, Padhraic Smyth, and Christopher KI Williams. Automating data science. *Communications of the ACM*, 2022.
  - Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. arXiv preprint arXiv:2411.15594, 2024.
  - Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. DS-agent: Automated data science by empowering large language models with case-based reasoning. *International Conference on Machine Learning*, 2024.
  - Md Mahadi Hassan, Alex Knipper, and Shubhra Kanti Karmaker Santu. Chatgpt as your personal data scientist. *arXiv preprint arXiv:2305.13657*, 2023.
  - Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Chenxing Wei, Danyang Li, Jiaqi Chen, Jiayi Zhang, et al. Data interpreter: An llm agent for data science. *arXiv* preprint arXiv:2402.18679, 2024.
  - Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing Xu, Ming Zhu, et al. Infiagent-dabench: Evaluating agents on data analysis tasks. *arXiv* preprint arXiv:2401.05507, 2024.
  - Yiming Huang, Jianwen Luo, Yan Yu, Yitong Zhang, Fangyu Lei, Yifan Wei, Shizhu He, Lifu Huang, Xiao Liu, Jun Zhao, et al. Da-code: Agent data science code generation benchmark for large language models. *arXiv preprint arXiv:2410.07331*, 2024.
  - Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. Aide: Ai-driven exploration in the space of code. *arXiv preprint arXiv:2502.13138*, 2025.
  - Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *International Conference on Learning Representations*, 2024.
  - Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. Dsbench: How far are data science agents to becoming data science experts? *International Conference on Learning Representations*, 2025.
  - Eugenie Lai, Gerardo Vitagliano, Ziyu Zhang, Sivaprasad Sudhir, Om Chabra, Anna Zeng, Anton A Zabreyko, Chenning Li, Ferdi Kossmann, Jialin Ding, et al. Kramabench: A benchmark for ai systems on data-to-insight pipelines over data lakes. *arXiv preprint arXiv:2506.06541*, 2025.
  - Jinhyuk Lee, Feiyang Chen, Sahil Dua, Daniel Cer, Madhuri Shanbhogue, Iftekhar Naim, Gustavo Hernández Ábrego, Zhe Li, Kaifeng Chen, Henrique Schechter Vera, et al. Gemini embedding: Generalizable embeddings from gemini. *arXiv preprint arXiv:2503.07891*, 2025.
  - Fei Li and Hosagrahar V Jagadish. Constructing an interactive natural language interface for relational databases. *Proceedings of the VLDB Endowment*, 2014.
- Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can Ilm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 2023.

- Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can Ilm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 2024.
  - Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 2022.
  - Karime Maamari, Fadhil Abubaker, Daniel Jaroslawicz, and Amine Mhedhbi. The death of schema linking? text-to-sql in the age of well-reasoned language models. *arXiv preprint arXiv:2408.07702*, 2024.
  - Friso Kingma Martin Iglesias, Alex Egg. Data agent benchmark for multi-step reasoning (DABStep), 2025.
  - Jaehyun Nam, Jinsung Yoon, Jiefeng Chen, Jinwoo Shin, Sercan Ö Arık, and Tomas Pfister. Mlestar: Machine learning engineering agent via search and targeted refinement. *arXiv preprint arXiv:2506.15692*, 2025.
  - Zhijie Nie, Zhangchi Feng, Mingxin Li, Cunwang Zhang, Yanzhao Zhang, Dingkun Long, and Richong Zhang. When text embedding meets large language model: a comprehensive survey. *arXiv* preprint arXiv:2412.09165, 2024.
  - Shuyin Ouyang, Dong Huang, Jingwen Guo, Zeyu Sun, Qihao Zhu, and Jie M Zhang. Ds-bench: A realistic benchmark for data science code generation. *arXiv preprint arXiv:2505.15621*, 2025.
  - Mohammadreza Pourreza and Davood Rafiei. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *Advances in Neural Information Processing Systems*, 2023.
  - Mohammadreza Pourreza, Hailong Li, Ruoxi Sun, Yeounoh Chung, Shayan Talaei, Gaurav Tarlok Kakkar, Yu Gan, Amin Saberi, Fatma Ozcan, and Sercan O Arik. Chase-sql: Multi-path reasoning and preference optimized candidate selection in text-to-sql. *arXiv preprint arXiv:2410.01943*, 2024.
  - Iqbal H Sarker. Data science and analytics: an overview from data-driven smart computing, decision-making and applications perspective. *SN Computer Science*, 2021.
  - Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 2023.
  - Zhaohao Sun, Lizhe Sun, and Kenneth Strang. Big data analytics services for enhancing business intelligence. *Journal of Computer Information Systems*, 2018.
  - Shayan Talaei, Mohammadreza Pourreza, Yu-Chen Chang, Azalia Mirhoseini, and Amin Saberi. Chess: Contextual harnessing for efficient sql synthesis. *arXiv preprint arXiv:2405.16755*, 2024.
  - Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv* preprint arXiv: Arxiv-2305.16291, 2023.
  - Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *International Conference on Learning Representations*, 2024.
  - Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023.
  - Xu Yang, Xiao Yang, Shikai Fang, Bowen Xian, Yuante Li, Jian Wang, Minrui Xu, Haoran Pan, Xinpeng Hong, Weiqing Liu, et al. R&d-agent: Automating data-driven ai solution building through llm-powered automated research, development, and evolution. *arXiv preprint arXiv:2505.14738*, 2025.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. International Conference on Learning Representations, 2023. Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. arXiv preprint arXiv:1809.08887, 2018. Wenqi Zhang, Yongliang Shen, Weiming Lu, and Yueting Zhuang. Data-copilot: Bridging billions of data and humans with autonomous workflow. arXiv preprint arXiv:2306.07209, 2023. Yuge Zhang, Qiyang Jiang, Xingyu Han, Nan Chen, Yuqing Yang, and Kan Ren. Benchmarking data science agents. arXiv preprint arXiv:2402.17168, 2024. Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in neural information processing systems, 2023. 

# **ALGORITHMS**

648

649 650

696 697 698

700 701

# **Algorithm 1** DS-STAR

```
651
             1: Input: query q, data files \mathcal{D}, maximum number of refinement round M
652
             2: # Analyzing data files
653
             3: for i = 1 to N do
654
                       s_{\mathtt{desc}}^i = \mathcal{A}_{\mathtt{analyzer}}(\mathcal{D}_i)
             4:
655
                       d_i = \mathrm{exec}(s^i_{\mathtt{desc}})
             5:
656
             6: end for
657
             7: # Iterative plan generation and verification
658
             8: p_0 = \mathcal{A}_{planner}(q, \{d_i\}_{i=1}^N)
659
             9: s_0 = \mathcal{A}_{coder}(p_0, \{d_i\}_{i=1}^N)
660
            10: r_0 = \exp(s_0)
661
            11: p = \{p_0\}
662
            12: for k = 0 to M - 1 do
663
                       v = \mathcal{A}_{\mathtt{verifier}}(p, q, s_k, r_k)
664
                       \quad \textbf{if } v = \texttt{sufficient then} \\
            14:
665
            15:
                           break
                       else if v = insufficient then
666
            16:
                            w = \mathcal{A}_{\texttt{router}}(p, q, r_k, \{d_i\}_{i=1}^N)
            17:
667
                           if w \in \{0, \cdots, \operatorname{len}(p) - 1\} then
            18:
668
                                 l = w - 1
            19:
669
            20:
                            else if w = Add Step then
670
                                 l = k
            21:
671
                            end if
            22:
672
            23:
                            p \leftarrow \{p_0, \cdots, p_l\}
673
                           p_{l+1} = \mathcal{A}_{\texttt{planner}}(p, q, r_k, \{d_i\}_{i=1}^N)
            24:
674
            25:
                            p \leftarrow \{p_0, \cdots, p_l, p_{l+1}\}
675
                            s_{k+1} = \mathcal{A}_{\texttt{coder}}(p, q, s_k, \{d_i\}_{i=1}^N)
            26:
676
            27:
                            r_{k+1} = \mathsf{exec}(s_{k+1})
            28:
                      end if
677
            29: end for
678
            30: Output: Final solution s
679
680
```

# B BENCHMARK

Table 5: **Statistics of benchmark.** For the DA-Code benchmark (Huang et al., 2024), we report the average of number of data files required for each tasks in each domain.

Benchmark	Domain	# Tasks	# Hard Tasks	# Data files
DABStep (Martin Iglesias, 2025)		450	378	7
	Archeology	12	6	5
	Astronomy	12	6	1556
V	Biomedical	9	6	7
KramaBench (Lai et al., 2025)	Environment	20	14	37
	Legal	30	16	136
	Wildfire	21	15	23
	Data Insight	79	5	3.7
	Data Manipulation	73	23	5.8
DA Codo (Huang et al. 2024)	Data Wrangling	100	16	6.2
DA-Code (Huang et al., 2024)	Machine Learning	100	34	3.9
	Statistical Analysis	70	10	4.5
	Visualization	78	14	4.5

We evaluate DS-STAR against several alternatives using two challenging benchmarks: DABStep (Martin Iglesias, 2025) and KramaBench (Lai et al., 2025).

- **DABStep** is composed of 450 tasks (72 easy and 378 hard) that require analyzing a shared set of seven data files. The ground-truth labels for these tasks are hidden, and evaluation is performed by submitting results to an official server, ensuring a blind assessment.
- **KramaBench** tests an agent's ability to perform autonomous data discovery. It contains tasks across six distinct domains, where each domain includes up to 1,556 data files. This setup requires the agent to automatically identify and select the relevant data files for given task.
- **DA-Code** is structured into three main categories: data wrangling, machine learning, and exploratory data analysis (EDA), with the EDA category further divided into data manipulation, data insights, visualization, and statistical analysis. See Appendix E for the example tasks.

Table 5 provides detailed statistics for both benchmarks.

#### C COST ANALYSIS

Table 6: Cost analysis on DABStep. We report the average cost when using Gemini-2.5-Pro.

Method	# LLM calls	# Input tokens	# Output tokens	Cost (\$)
ReAct (Yao et al., 2023)	7.1	44691	2928	0.09
DA-Agent (Huang et al., 2024)	8.8	39999	4123	0.09
DS-STAR (Ours)	12.7	154669	3373	0.23

In this section, we quantitatively analyze the LLM usage cost. We evaluate DS-STAR on 10 development tasks from the DABStep benchmark (Martin Iglesias, 2025), reporting the average number of LLM calls, input tokens, and output tokens.

We acknowledge that DS-STAR incurs higher cost due to its increased token usage. As shown in Table 6, DS-STAR requires 3.5 times more input tokens than the ReAct baseline. This increase is primarily because our method utilizes comprehensive analytic descriptions of each data file, which inherently necessitates more input tokens. Nevertheless, this translates to a modest practical cost; using Gemini-2.5-Pro (priced at \$1.25/\$10.00 per 1M input/output tokens³), the expense for DS-STAR averages only \$0.23 per data science task. We argue this higher cost is a justifiable trade-off, as DS-STAR consistently and significantly achieves better results than ReAct.

<sup>&</sup>lt;sup>3</sup>https://ai.google.dev/gemini-api/docs/pricing

757 758

759 760

761

762

763

764

765

766

767

768 769

770 771

772

773

774

775

776

777

778

779

780

781

783 784

786

788

789 790

791

792

793

794

795

796 797

798

800

801

802

803

804

805

806

807

808

809

# D QUALITATIVE RESULTS ON THE ANALYZER AGENT

#### D.1 CODE GENERATED BY THE ANALYZER AGENT

```
import ison
def describe_json_data(file_path):
    Loads a JSON file, prints its full content, and then describes its structure
    and essential information, including column names (keys) and example data.
     # Load the JSON data from the file
    # Using utf-8 encoding is a good practice for JSON.
with open(file_path, 'r', encoding='utf-8') as f:
         data = json.load(f)
    print(f"--- Describing content of: {file_path} ---")
    print("\n")
     # 1. Print the full loaded content as per requirement
    # "The Python code should print out the content of data/merchant_data.json."
print("--- Full Loaded Content ---")
     ^{	ilde{+}} Using json.dumps for a pretty-printed representation of the loaded Python object.
    # ensure_ascii=False allows for proper display of unicode characters.
print(json.dumps(data, indent=4, ensure_ascii=False))
    print("-" * 50)
    print("\n")
    # 2. Print Essential Information
print("--- Essential Information ---")
    data_type = type(data)
    print(f"Overall Data Type: {data_type.__name__}}")
    if isinstance(data, dict):
         print("The data is a single JSON object (dictionary).")
         keys = list(data.keys())
         print(f"Number of top-level keys: {len(keys)}")
             print("The dictionary is empty.")
         else:
              # Requirement: "print out all the column names."
             print("Top-level keys (can be considered as 'column names' for the main object):") for key in keys:
             \# Requirement: "If there are too many structured data, print out just few examples." print("\n--- Example Key-Value Pairs (first 5 if more exist) ---")
              count = 0
              for key, value in data.items():
                  if count < 5:</pre>
                       value str = str(value)
                         Truncate long values for a concise preview
                      preview = (value_str[:70] + '...') if len(value_str) > 70 else value_str
print(f" {key}: {preview}")
                       count += 1
                  else:
                      print("
                                 ... (more key-value pairs exist)")
             if count == 0 and not keys: # Should be caught by 'if not keys:' earlier
                   print(" (Dictionary is empty - no key-value pairs to show)")
    elif isinstance(data, list):
         print("The data is a JSON array (list).")
num_items = len(data)
         print(f"Number of items in the list: {num items}")
         if num_items == 0:
             print("The list is empty.")
              first item = data[0]
             print(f"Type of the first item in the list: {type(first item). name }")
              # Check if it's a list of dictionaries (common for structured data records)
              if isinstance(first_item, dict):
                  print("The list appears to contain JSON objects (dictionaries), suggesting structured
                  ⇔ data.")
                  # Requirement: "print out all the column names."
# For a list of objects, column names are typically the keys of these objects.
                   # We'll use the keys from the first object as representative column names.
                  first_item_keys = list(first_item.keys())
                  if not first_item_keys:
                       print("The first object in the list is an empty dictionary (no column names to infer).")
                       print("Keys of the first object (assumed to be common 'column names' for the items):")
                       for key in first_item_keys:
    print(f" - {key}")
                   # Requirement: "If there are too many structured data, print out just few examples."
```

```
810
                               num_examples_to_show = min(3, num_items) # Show up to 3 examples
811
                                               Example Items (first {num_examples_to_show} of {num_items}) ---")
                               for i in range(num_examples_to_show):
812
                                   print(f"Item {i+1}:")
                                    # Pretty print each example item. Truncate if an individual item is very large.
example_item_str = json.dumps(data[i], indent=2, ensure_ascii=False)
813
                                    lines = example_item_str.split('\n')
814
                                   if len(lines) > 15: \# Arbitrary limit: max 15 lines per example item summary print('\n'.join(lines[:15]))
815
                                        print("
                                                     ... (item content truncated for brevity in this summary)")
816
                                    else:
                                       print(example_item_str)
817
818
                               if num_items > num_examples_to_show:
    print(f"\n... and {num_items - num_examples_to_show} more items in the list.")
819
                          # Else, it's a list of other types (e.g., strings, numbers - unstructured or semi-structured)
820
821
                               print("The list contains non-object items (e.g., strings, numbers, booleans, or mixed
                               822
                               print("This is often considered unstructured or semi-structured data.")
                               num_examples_to_show = min(5, num_items) # Show up to 5 examples
print(f"\n--- Example Items (first {num_examples_to_show}) of {num_items}) ---
823
                               for i in range(num_examples_to_show):
824
                                    item str = str(data[i])
                                    # Truncate long string representations for a concise preview
825
                                    preview = (item_str[:70] + '...') if len(item_str) > 70 else item_str
826
                                    print(f"Item {i+1}: {preview}")
                               if num_items > num_examples_to_show:
    print(f"\n... and {num_items - num_examples_to_show} more items in the list.")
828
829
                 elif isinstance(data, (str, int, float, bool)) or data is None:
                      # For simple scalar types, the "Full Loaded Content" print is the primary information.
830
                      # This section just confirms the nature of the data.
                      print \ (f" \ The \ data \ is \ a \ simple \ JSON \ scalar \ value \ (not \ an \ object \ or \ array) \ .") 
831
                      if data is None:
832
                          print("Specifically, the JSON content is 'null'.")
                      \ensuremath{\sharp} The actual value was already printed in the "Full Loaded Content" section.
833
834
                     # This case should ideally not be reached if json.load() was successful
# and returned a standard Python type corresponding to a valid JSON root element.
835
                     print("The data type is unusual for a root JSON element (e.g., a custom object if json.load was \hookrightarrow hooked).")
836
                     # The actual value was already printed in the "Full Loaded Content" section.
837
                 print("\n" + "-" * 50)
838
                 print("--- End of Description ---")
839
                 _name__ == "__main__
840
                   The path to the JSON file.
                # This script assumes 'data/merchant_data.json' exists relative to where the script is run.
841
                 # For example, if the script is in /home/user/scripts/, it expects
                    /home/user/scripts/data/merchant_data.json
                 json_file_path = 'data/merchant_data.json
843
                 # Call the function to load and describe the JSON data.
844
                 # As per requirements, no try-except blocks are used for error handling.
                       the file doesn't exist, a FileNotFoundError will be raised.
845
                 # If the file's content is not valid JSON, a json.JSONDecodeError will be raised. # These errors will halt the program, and debugging can proceed from there.
846
                 describe_json_data(json_file_path)
847
```

Listing 1: Python code for generating the description of 'merchant\_data.json' in DABStep.

```
864
865
            import pandas as pd
            import numpy as np
866
            import io
867
            import os
868
            # --- Configuration --
            # Define the path to the Excel file to be analyzed.
869
            \# This script assumes the file is located at {\rm \dot{'}data/1-s2.0-S0092867420301070-mmc4.xlsx'.}
            file_path = 'data/1-s2.0-S0092867420301070-mmc4.xlsx'
870
871
            # --- Analyst's Report Generation ---
872
            # Print a formatted header for the analysis report.
873
            print(f" Expert Data Analysis Report for: {file_path}")
874
875
            # Check for file existence to provide a clear startup message.
            # The script will still raise an error if not found, as per requirements.
if not os.path.exists(file_path):
876
                print(f"\n[STATUS] File not found at the specified path: '{file_path}'")
print("Please ensure the file exists and the path is correct before execution.")
# The script will exit here due to a FileNotFoundError from pandas,
877
878
                 # allowing for debugging as requested.
879
            # --- 1. Loading Data and Sheet Information --
880
881
            print("\n--- 1. Loading & Sheet Information ---")
            # Use pandas.ExcelFile to inspect the workbook without loading all data initially.
882
            \# This is efficient for understanding the structure of a multi-sheet Excel file.
            excel_file = pd.ExcelFile(file_path)
sheet_names = excel_file.sheet_names
883
884
            print(f"Found {len(sheet_names)} sheet(s) in the workbook: {sheet_names}")
885
            # For this analysis, we will focus on the first sheet.
            first_sheet_name = sheet_names[0]
print(f"Loading data from the first sheet: '{first_sheet_name}'")
886
            df = pd.read_excel(excel_file, sheet_name=first_sheet_name)
887
            print("Data loaded successfully.")
888
            # --- 2. High-Level Overview ---
889
            print("\n--- 2. Data Dimensions and Structure ---")
890
            rows, cols = df.shape
            print(f"The dataset contains {rows} rows and {cols} columns.")
891
892
            # An analyst's initial assessment of data structure.
            if cols == 1:
893
                print("Observation: With only one column, the data might be unstructured (e.g., a list of texts or
            → IDs).")
elif 1 < cols < 5:
894
895
                 print("Observation: With a few columns, the data appears to be structured but simple.")
            else:
896
                print("Observation: With multiple columns, the data appears to be well-structured.")
897
            # --- 3. Column Analysis ---
898
            print("\n--- 3. Column Names ---")
899
            print("The following columns are present in the dataset:")
            # Using a list comprehension for clean printing of column names.
print(" | ".join([f"'(col)'" for col in df.columns]))
900
901
902
            # --- 4. Data Types and Non-Null Counts ---
903
            print("\n--- 4. Data Types & Memory Usage ---")
# Capture the output of df.info() to present it cleanly within the report.
904
            buffer = io.StringIO()
905
            df.info(buf=buffer)
            info str = buffer.getvalue()
906
            print (info str)
907
            # --- 5. Missing Value Analysis ---
908
            print("\n--- 5. Missing Values Summary ---")
            missing_values = df.isnull().sum()
909
            total_missing = missing_values.sum()
910
            if total_missing == 0:
    print("Excellent! No missing values were found in the dataset.")
911
                print(f"Found a total of {total_missing} missing values.")
# Display columns that have at least one missing value.
912
913
                 missing_per_column = missing_values[missing_values > 0].sort_values(ascending=False)
print("Columns with missing values:")
914
                 print(missing_per_column.to_string())
915
            # --- 6. Descriptive Statistics for Numerical Data ---
916
            print("\n--- 6. Descriptive Statistics (for Numerical Columns) ---")
917
             # The describe() method provides a powerful summary of central tendency,
            # dispersion, and shape of the distribution for numerical columns.
```

```
# It automatically excludes non-numerical columns.
descriptive_stats = df.describe()
if descriptive_stats.empty:
    print("No numerical columns found to generate statistics for.")
    print (descriptive_stats)
# --- 7. Data Preview ---
print("\n--- 7. Data Sample (First 5 Rows) ---")
# Displaying the head of the DataFrame is crucial for getting a feel
# for the actual content and format of the data.
print (df.head())
# --- 8. Analysis of Categorical / Object Data ---
print("\n--- 8. Unique Values in Categorical Columns ---")
# Select columns with 'object' dtype, which are often categorical.
object_cols = df.select_dtypes(include=['object']).columns
if len(object_cols) == 0:
    print("No categorical (object type) columns found for value count analysis.")
else:
     print(f"Analyzing {len(object_cols)} categorical column(s): {list(object_cols)}\n")
     for col in object_cols:
    unique_count = df[col].nunique()
    print(f"-> Column '{col}' has {unique_count} unique values.")
           # To avoid flooding the output, only show value counts for columns with few unique values.
          if unique_count > 1 and unique_count <= 15:
    print(" Value Counts:")</pre>
               print(" Value Counts:")
print(df[col].value_counts().to_string())
print("-" * 20)
\# --- End of Report ---
                          End of Analysis Report")
print ("======="")
```

Listing 2: Python code for generating the description of '1-s2.0-S0092867420301070-mmc4.xlsx' in KramaBench.

1025

# D.2 EXAMPLES OF GENERATED DATA DESCRIPTIONS

```
973
974
       File: data/merchant_data.json (6857 bytes)
975
       Top-level type: array
976
       root: array with 30 items
       Element types: object=30
977
        Columns (5): ['account_type', 'acquirer', 'capture_delay', 'merchant',
978
        → 'merchant_category_code']
979
       Sample rows (first 5):
980
981
            "merchant": "Crossfit_Hanna",
982
            "capture_delay": "manual",
983
            "acquirer": [
984
              "gringotts",
985
              "the_savings_and_loan_bank",
986
              "bank_of_springfield",
              "dagoberts_vault"
987
988
            "merchant_category_code": 7997,
            "account_type": "F"
990
991
            "merchant": "Martinis_Fine_Steakhouse",
992
            "capture_delay": "immediate",
993
            "acquirer": [
994
              "dagoberts_geldpakhuis",
995
              "bank_of_springfield"
996
            "merchant_category_code": 5812,
997
            "account_type": "H"
998
          },
999
1000
            "merchant": "Belles_cookbook_store",
1001
            "capture_delay": "1",
            "acquirer": [
1002
              "lehman_brothers"
1003
            "merchant_category_code": 5942,
1005
            "account_type": "R"
1006
          },
1007
            "merchant": "Golfclub_Baron_Friso",
1008
            "capture_delay": "2",
1009
            "acquirer": [
1010
              "medici"
1011
            "merchant_category_code": 7993,
1012
            "account_type": "F"
1013
1014
1015
            "merchant": "Rafa_AI",
1016
            "capture_delay": "7",
            "acquirer": [
1017
              "tellsons_bank"
1018
1019
            "merchant_category_code": 7372,
1020
            "account_type": "D"
1021
1022
1023
1024
```

Listing 3: Description of 'merchant\_data.json' in DABstep.

```
1026
1027
       f'''
1028
       ______
1029
        Expert Data Analysis Report for:

→ data/1-s2.0-S0092867420301070-mmc4.xlsx

1030
       _____
1031
1032
       --- 1. Loading & Sheet Information ---
1033
       Found 4 sheet(s) in the workbook: ['README', 'A-Variants', 'B-Novel
       → Splice Junctions', 'C-Alternate Splice Junctions']
1034
       Loading data from the first sheet: 'README'
1035
       Data loaded successfully.
1036
1037
       --- 2. Data Dimensions and Structure ---
1038
       The dataset contains 25 rows and 2 columns.
       Observation: With a few columns, the data appears to be structured but
1039
       \hookrightarrow simple.
1040
1041
       --- 3. Column Names ---
1042
       The following columns are present in the dataset:
1043
       'IN THIS FILE: a list of variants and alternate splice junctions for
1044
       → which we found peptide evidence. For more information, please see
       → STAR Methods' | 'Unnamed: 1'
1045
1046
       --- 4. Data Types & Memory Usage ---
1047
       <class 'pandas.core.frame.DataFrame'>
1048
       RangeIndex: 25 entries, 0 to 24
       Data columns (total 2 columns):
1049
           Column
1050
       → Non-Null Count Dtype
1051
       ____
1052
1053
       Ω
           IN THIS FILE: a list of variants and alternate splice junctions for
          which we found peptide evidence. For more information, please see
1054
       \hookrightarrow
          STAR Methods 21 non-null
                                        object
1055
        1
            Unnamed: 1
1056

→ 21 non-null

                            object
1057
       dtypes: object(2)
1058
       memory usage: 528.0+ bytes
1059
1060
       --- 5. Missing Values Summary ---
1061
       Found a total of 8 missing values.
1062
       Columns with missing values:
1063
       IN THIS FILE: a list of variants and alternate splice junctions for
       \hookrightarrow which we found peptide evidence. For more information, please see
1064
       \hookrightarrow STAR Methods
1065
       Unnamed: 1
1066
       \hookrightarrow 4
1067
1068
       --- 6. Descriptive Statistics (for Numerical Columns) ---
              IN THIS FILE: a list of variants and alternate splice junctions
1069
          for which we found peptide evidence. For more information, please
1070
       \hookrightarrow see STAR Methods
                              Unnamed: 1
1071
       count.
                                                                2.1
1072

→ 21

1073
       unique
                                                                21
1074
       top
                                                             Sheet
1075
       \hookrightarrow Description
1076
       freq
                                                                  1
1077
       \hookrightarrow 1
1078
       --- 7. Data Sample (First 5 Rows) ---
1079
```

```
1080
        IN THIS FILE: a list of variants and alternate splice junctions for
1081
       \,\hookrightarrow\, which we found peptide evidence. For more information, please see
1082
           STAR Methods
                                                                    Unnamed: 1
                                                           NaN
1083
           NaN
1084
       1
                                                           NaN
1085
           NaN
       \hookrightarrow
1086
       2
                                                         Sheet
1087
           Description
1088
       3
                                                    A-Variants
           Both single amino acid variants (SAAVs) and in...
       \hookrightarrow
1089
       4
                                    B-Novel Splice Junctions
1090
       → Splice junctions where neither boundary was at...
1091
1092
       --- 8. Unique Values in Categorical Columns ---
       Analyzing 2 categorical column(s): ['IN THIS FILE: a list of variants
1093
       \,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\, and alternate splice junctions for which we found peptide evidence.
1094
       \hookrightarrow For more information, please see STAR Methods', 'Unnamed: 1']
1095
1096
       -> Column 'IN THIS FILE: a list of variants and alternate splice
1097
       → junctions for which we found peptide evidence. For more information,
       → please see STAR Methods' has 21 unique values.
1098
       -> Column 'Unnamed: 1' has 21 unique values.
1099
1100
1101
                     End of Analysis Report
1102
       _____
1103
```

Listing 4: Description of '1-s2.0-S0092867420301070-mmc4.xlsx' in KramaBench.

# E QUALITATIVE RESULTS FOR EACH TASK TYPE

This section provides qualitative results. It includes queries on DA-Code benchmark, solution code generated by DS-STAR, output generated by DS-STAR, and ground-truth answer. Specifically, it presents examples for each task, such as statistical analysis, data insights, data wrangling, machine learning, and visualization tasks.

#### E.1 STATISTICAL ANALYSIS

1134

1135 1136

1137

1138

1139 1140

1141

1142

1143

1144

1145 1146

1147 1148

1186

```
f'''Could you analyze our Business Case dataset to perform a hypothesis

→ test comparing bike rentals across different weather conditions? I

→ need the results formatted and saved in a file named 'weather.csv'

→ following the requirement in tips.md.'''
```

Listing 5: User's question.

```
import pandas as pd
1149
            from scipy.stats import kruskal
1150
1151
            # Define the file path for the input dataset
           input_file_path = 'data/yulu_bike_sharing_dataset.csv'
1152
1153
            # Load the dataset into a pandas DataFrame
           yulu_df = pd.read_csv(input_file_path)
1154
            # The goal is to perform a Kruskal-Wallis H-test to compare the distribution of bike rentals ('count')
1155
           # across different weather conditions ('weather').
1156
             The 'weather' column contains numerical codes for different weather conditions.
1157
           # First, we need to find the unique weather conditions in the dataset.
unique_weather_conditions = sorted(yulu_df['weather'].unique())
1158
           \# Next, we create a list of data samples, where each sample is an array of 'count' \# values corresponding to a specific weather condition.
1159
1160
           weather_groups = [yulu_df['count'][yulu_df['weather'] == weather] for weather in unique_weather_conditions]
1161
            \# Perform the Kruskal-Wallis H-test. The \star operator unpacks the list of groups
           # into separate arguments for the function.
statistic, p_value = kruskal(*weather_groups)
1162
1163
           # Define the significance level (alpha) as specified in the guidelines.
1164
           alpha = 0.05
1165
           # Based on the p-value, determine the hypothesis conclusion and the comment
             as per the format requirements in 'data/tips.md'.
1166
           if p_value < alpha:</pre>
                hypothesis_conclusion = "Strong Evidence to Reject Null Hypothesis"
1167
                comment_on_results = "Different distributions
1168
               hypothesis_conclusion = "No Evidence to Reject Null Hypothesis"
1169
                comment_on_results = "No evidence of different distributions
1170
           # Create a dictionary containing the results in the required format.
1171
                'Test Type': ["Non Parametric"],
'Kruskal Statistic': [statistic],
1172
                'P-Value': [p_value],
1173
                'Hypothesis': [hypothesis_conclusion],
                'Comment': [comment_on_results]
1174
1175
            # Create a new pandas DataFrame from the results dictionary.
1176
           weather_results_df = pd.DataFrame(results_data)
1177
           # Define the output directory and file path.
output_dir = 'final'
1178
           output_file_path = os.path.join(output_dir, 'weather.csv')
1179
            # Ensure the output directory exists.
1180
           os.makedirs(output_dir, exist_ok=True)
1181
            # Save the resulting DataFrame to the specified CSV file,
1182
            # ensuring the index is not written to the file.
           weather_results_df.to_csv(output_file_path, index=False)
1183
1184
            # print(f"Analysis complete. Results saved to {output_file_path}")
           # print(weather_results_df.to_string())
1185
```

Listing 6: Solution code.

Table 7: Generated output by DS-STAR in csv format.

Test Type	Kruskal Statistic	P-Value	Hypothesis	Comment
Non Parametric	205.00216514479087	3.501611300708679e-44	Strong Evidence to Reject Null Hypothesis	Different distributions

#### E.2 DATA INSIGHT

```
f'''Could you please analyze the matches between Argentina and Brazil?

→ Record the total number of matches, the number of wins for each

→ team, and the number of draws. Indicate if Brazil has won more

→ matches ("Yes" or "No"). Present your findings in the following

→ format:

"""

{

"Total Matches between Argentina and Brazil": [...],

"Arg Wins": [...],

"Brazil Wins": [...],

"Draws": [...],

"Result": [...]

}

""""'''
```

Listing 7: User's question.

```
1210
1211
            import pandas as pd
           import numpy as np
import json
1212
           import os
1213
           # Ensure the 'final' directory exists for the output file
output_dir = 'final'
1214
1215
           os.makedirs(output dir, exist ok=True)
1216
           # Load the results.csv file into a pandas DataFrame
           results df = pd.read csv('data/results.csv')
1217
1218
           # Filter for matches where the home team is Argentina and the away team is Brazil
           arg_vs_bra = results_df[
1219
                (results_df['home_team'] == 'Argentina') & (results_df['away_team'] == 'Brazil')
1220
           # Filter for matches where the home team is Brazil and the away team is Argentina
bra_vs_arg = results_df[
1221
1222
               (results_df['home_team'] == 'Brazil') & (results_df['away_team'] == 'Argentina')
1223
           # Concatenate the two DataFrames to get all matches between Argentina and Brazil
1224
           argentina_brazil_matches = pd.concat([arg_vs_bra, bra_vs_arg])
1225
           # Sort the combined DataFrame by date for a chronological view
1226
           argentina_brazil_matches = argentina_brazil_matches.sort_values(by='date').copy()
1227
           # Define the conditions for determining the winner based on scores
1228
           # A draw in regular/extra time is considered a draw, regardless of a penalty shootout outcome.
1229
               argentina_brazil_matches['home_score'] > argentina_brazil_matches['away_score'],
argentina_brazil_matches['away_score'] > argentina_brazil_matches['home_score']
1230
1231
           # Define the corresponding outcomes (the winning team's name)
1232
               argentina brazil matches['home team'],
1233
               argentina_brazil_matches['away_team']
1234
1235
           # Create the 'winner' column using np.select, with 'Draw' as the default
           argentina_brazil_matches['winner'] = np.select(conditions, choices, default='Draw')
1236
             Count the total wins for each team and the number of draws
1237
           win_counts = argentina_brazil_matches['winner'].value_counts()
1238
            # Calculate the total number of matches, wins for each team, and draws from the win_counts Series
1239
           total_matches = int(win_counts.sum())
argentina_wins = int(win_counts.get('Argentina', 0))
1240
           brazil_wins = int(win_counts.get('Brazil', 0))
1241
           draws = int(win_counts.get('Draw', 0))
           # Determine the 'Result' value. 'Yes' if Brazil has more wins, otherwise 'No'.
```

```
1242
           result_value = 'Yes' if brazil_wins > argentina_wins else 'No'
1243
           # Construct the dictionary to store the analysis results with the specified keys
1244
           analysis_results = {
1245
               "Total Matches between Argentina and Brazil": total_matches,
               "Arg Wins": argentina_wins,
1246
               "Brazil Wins": brazil_wins,
               "Draws": draws,
"Result": result_value
1247
1248
1249
           # Define the output file path
          output_path = os.path.join(output_dir, 'result.json')
1250
           # Save the final results dictionary to a JSON file in the 'final' directory
1251
           with open(output_path, 'w') as f:
               ison.dump(analysis results, f, indent=4)
1252
```

#### Listing 8: Solution code.

```
f'''{
    "Total Matches between Argentina and Brazil": [109],
    "Arg Wins": [40],
    "Brazil Wins": [43],
    "Draws": [26],
    "Result": ["Yes"]
}'''
```

Listing 9: **Generated output** by DS-STAR, in json format.

#### E.3 DATA MANIPULATION

1253

1258

1259

1260

1261

1262 1263

1264

12651266

1267

1269

1270 1271

```
f'''Calculate the total quantity sold and total sales revenue for each

→ bike category from 2016 to 2018. Ensure that total sales revenue

→ accounts for any discounts applied to each order. Write the results

→ into result.csv, following the exact structure and formatting of

→ sample_result.csv.'''
```

Listing 10: User's question.

```
1273
            import pandas as pd
import os
1274
1275
            # Ensure the output directory exists
1276
           os.makedirs('final', exist_ok=True)
1277
            # Load the required data files into pandas DataFrames
           order_items_df = pd.read_csv('data/order_items.csv')
1278
           products_df = pd.read_csv('data/products.csv')
categories_df = pd.read_csv('data/categories.csv')
1279
            sample_result_df = pd.read_csv('data/sample_result.csv')
1280
            # Calculate the sales revenue for each item, accounting for discounts.
1281
           # sales_revenue = quantity * list_price * (1 - discount)
order_items_df['sales_revenue'] = order_items_df['quantity'] * order_items_df['list_price'] * (1 -
1282
            → order_items_df['discount'])
1283
            # Merge the order items with product information to get model year and category id
1284
           merged_df = pd.merge(order_items_df, products_df, on='product_id')
1285
            # Merge with category information to get category name
1286
           merged_df = pd.merge(merged_df, categories_df, on='category_id')
1287
             Filter the DataFrame to include only products with model years from 2016 to 2018
           filtered_df = merged_df[(merged_df['model_year'] >= 2016) & (merged_df['model_year'] <= 2018)]</pre>
1288
1289
            # Group the filtered data by 'category_name' and aggregate the total quantity and total sales revenue.
           category_summary = filtered_df.groupby('category_name').agg(
    total_quantity=('quantity', 'sum'),
1290
                total_price=('sales_revenue', 'sum')
1291
           ).reset_index()
1292
             Use the sample_result.csv to define the structure and order of the final output.
1293
           \# Perform a left merge to ensure all categories from the sample are present in the final result, in the \hookrightarrow correct order.
1294
           final_df = pd.merge(sample_result_df[['category_name']], category_summary, on='category_name', how='left')
1295
            # Fill any potential NaN values with 0 for categories that had no sales in the filtered period.
           final_df[['total_quantity', 'total_price']] = final_df[['total_quantity', 'total_price']].fillna(0)
```

```
1296
1297
            # Ensure the column order matches the sample_result.csv file: ['total_quantity', 'total_price',
                 'category name
1298
            final_df = final_df[['total_quantity', 'total_price', 'category_name']]
1299
            # Convert the 'total_quantity' column to integer type as it represents a count of items.
final_df['total_quantity'] = final_df['total_quantity'].astype(int)
1300
            # Define the output path
            output_path = 'final/result.csv'
1302
1303
            # Save the final DataFrame to the specified path, without the index.
            final_df.to_csv(output_path, index=False)
1304
            print(f"Successfully generated '{output_path}'")
1305
            print("Content of the generated file:")
print(final_df)
1306
```

Listing 11: **Solution code.** 

Table 8: **Generated output** by DS-STAR in csv format.

total_quantity	total_price	category_name
1755	2715079.53	Mountain Bikes
559	1665098.49	Road Bikes
2063	995032.62	Cruisers Bicycles
315	916684.78	Electric Bikes
394	711011.84	Cyclocross Bicycles
813	394020.10	Comfort Bicycles
1179	292189.20	Children Bicycles

#### E.4 DATA WRANGLING

```
f'''Manage missing data in the NYC_Open_Data_Parking_Violations dataset

by either dropping or replacing it. Drop rows if necessary, and drop

columns only if most entries are empty. Replace missing values with

the mean, mode, or other appropriate methods. Save the cleaned data

to 'cleaned_data.csv'.'''
```

Listing 12: User's question.

```
1331
              import pandas as pd
1332
              import numpy as np
             import os
1333
              # The user prompt mentions cleaning the 'NYC_Open_Data_Parking_Violations' dataset,
1334
             # however, the provided data file is 'data/data.csv', which contains automobile data.
# This script will proceed with cleaning the provided 'data/data.csv' file based on the
1335
             # general data cleaning guidelines in the prompt.
1336
1337
             input_file = 'data/data.csv'
output_dir = 'final'
1338
             output_file = os.path.join(output_dir, 'cleaned_data.csv')
              # Create the output directory if it does not exist
1340
             if not os.path.exists(output_dir):
                  os.makedirs(output_dir)
1341
             # Load the data, interpreting '?' as Not a Number (NaN)
df = pd.read_csv(input_file, na_values='?')
1342
1343
              # --- Data Cleaning Steps ---
1344
             # Guideline: "Drop rows if necessary"
# The 'price' column is crucial for most analyses. Rows with a missing price are dropped.
df.dropna(subset=['price'], inplace=True)
1345
1346
1347
              # Guideline: "Replace missing values with the mean, mode, or other appropriate methods"
1348
              # Identify columns for mean and mode imputation
             # Numerical columns with missing values to be imputed with the mean
numeric_cols_to_impute_mean = ['normalized-losses', 'bore', 'stroke', 'horsepower', 'peak-rpm']
1349
```

```
1350
           \# Categorical columns with missing values to be imputed with the mode
1351
           categorical_cols_to_impute_mode = ['num-of-doors']
1352
             Impute numerical columns with the mean
1353
            for col in numeric_cols_to_impute_mean:
                # Convert column to numeric type. This is necessary because the '?' # values caused pandas to read them as 'object' type.
1354
                df[col] = pd.to_numeric(df[col])
# Calculate the mean of the column (NaNs are automatically excluded)
1355
                 mean_val = df[col].mean()
1356
                # Fill the missing (NaN) values with the calculated mean
df[col].fillna(mean_val, inplace=True)
1357
1358
            # Impute categorical columns with the mode
            for col in categorical_cols_to_impute_mode:
1359
                # Calculate the mode (most frequent value) of the column.
# .mode() returns a Series, so we take the first element.
1360
                mode_val = df[col].mode()[0]
1361
                 # Fill the missing (NaN) values with the calculated mode
                df[col].fillna(mode_val, inplace=True)
1362
            # Guideline: After handling missing values, ensure data types are correct.
# The 'price' column was loaded as 'object' and needs to be converted to a numeric type.
1363
1364
           df['price'] = df['price'].astype(float)
1365
            # --- Save the Cleaned Data ---
1366
1367
            # Save the final cleaned DataFrame to 'final/cleaned_data.csv'
            \sp{\#} The index is not saved as it's not part of the original data.
1368
           df.to_csv(output_file, index=False)
1369
            # --- Verification ---
           # This part confirms that the cleaning was successful and the file was created.
1370
           print(f"Successfully cleaned the dataset and saved it to '{output_file}'"
1371
            # Load the cleaned data to verify its integrity
1372
            cleaned_df = pd.read_csv(output_file)
1373
           print("\n--- Verification of Cleaned Data ---")
1374
1375
            # Check for any remaining missing values
           missing_values_summary = cleaned_df.isnull().sum()
           total_missing = missing_values_summary.sum()
1377
            print("\nMissing values per column in the cleaned file:")
1378
           print (missing_values_summary)
           print(f"\nTotal remaining missing values: {total_missing}")
1379
            if total missing == 0:
1380
                print("\nVerification successful: The cleaned dataset has no missing values.")
1381
                print("\nVerification failed: The cleaned dataset still contains missing values.")
1382
```

Listing 13: Solution code.

#### Table 9: Generated output by DS-STAR in csv format.

```
ymbing semigracians with the control of the control
```

#### E.5 MACHINE LEARNING

1383

1388

1389 1390 1391

1392

1394

1395 1396

```
f'''As a contestant in the Bank customer data churn prediction

→ competition, you need to predict the data in test.csv according to

→ the competition requirements. Write the results into submission.csv

→ following the format specified in sample_submission.csv.'''
```

Listing 14: User's question.

```
import pandas as pd
import numpy as np
from sklearn.model_selection import StratifiedKFold
import lightgbm as lgb
from sklearn.metrics import roc_auc_score
import optuna
import os

# Create the 'final' directory if it doesn't already exist to store the output.
```

```
1404
                   os.makedirs('final', exist_ok=True)
1405
                   # Suppress Optuna's logging to keep the output clean during the study.
1406
                   optuna.logging.set_verbosity(optuna.logging.WARNING)
1407
                   # 1. Load All Data
1408
                   # Load data from the 'data/' directory.
                   original_train_df = pd.read_csv('data/train.csv')
test_df = pd.read_csv('data/test.csv')
1409
                   additional_data_df = pd.read_csv('data/Churn_Modelling.csv')
1410
1411
                   # The provided test.csv uses 'CustomerId' as the unique identifier.
                   # We will use 'CustomerId' for the submission file's 'id' column as required by the competition format.
1412
                   test_ids = test_df['CustomerId']
1413
1414
                   # 2. Preprocess Data: Augment and Clean
                   # Select a consistent set of feature columns and the target variable ('Exited') to ensure uniformity across
1415
                    → datasets.
                   -- dataset.

| dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | dataset. | datase
1416
1417
                   # Augment the training data by combining the original training set with the additional 'Churn_Modelling.csv'
1418
                   train df = pd.concat([original train df[feature cols], additional data df[feature cols]], ignore index=True)
1419
1420
                   # --- START: Feature Engineering ---
1421
                   def feature_engineer(df):
                           """Creates new features based on domain knowledge and interactions between existing features."""
1422
                           \sharp A small epsilon is added to denominators to prevent division by zero errors.
1423
                          df['BalanceSalaryRatio'] = df['Balance'] / (df['EstimatedSalary'] + epsilon)
                         df('BalanceSalaryRatio') = df('Balance') / (df('EstimatedSalary') + epsi
df('TenureByAge') = df('Tenure') / (df('Age') + epsilon)
df('CreditScoreGivenAge') = df('CreditScore') / (df('Age') + epsilon)
df('Loyalty') = df('Tenure') * df('IsActiveMember')
df('ProductsPerTenure') = df('NumOfProducts') / (df('Tenure') + epsilon)
1424
1425
1426
                          return df
1427
                   print("Applying feature engineering to training and test data...")
                   train_df = feature_engineer(train_df)
test_df = feature_engineer(test_df)
# --- END: Feature Engineering ---
1428
1429
1430
                   # Drop identifier columns ('CustomerId', 'Surname', 'id') that are not useful for the model.
# 'errors='ignore'' handles cases where a column might not exist in a dataframe.
drop_cols = ['CustomerId', 'Surname', 'id']
1431
1432
                   train_df = train_df.drop(columns=[col for col in drop_cols if col in train_df.columns and col != 'Exited'],
1433

→ errors='ignore')
                   test df = test df.drop(columns=[col for col in drop cols if col in test df.columns], errors='ignore')
1434
1435
                   # One-Hot Encode categorical features ('Geography', 'Gender') to convert them into a machine-readable
1436
                         format
                   categorical_features = ['Geography', 'Gender']
train_df = pd.get_dummies(train_df, columns=categorical_features, drop_first=True)
1437
                   test_df = pd.get_dummies(test_df, columns=categorical_features, drop_first=True)
1438
1439
                   # Separate the target variable ('Exited') from the training features.
                   train_labels = train_df['Exited']
train_df = train_df.drop(columns=['Exited'])
1440
1441
                   # Align the training and test dataframes to ensure they have the exact same feature columns in the same
1442
                   \hookrightarrow order. \# 'inner' join keeps only columns that are present in both dataframes.
1443
                   train_df, test_df = train_df.align(test_df, join='inner', axis=1)
1444
                   \mbox{\#} Define the final feature set (X) and target (y).
1445
                   X = train df
                   y = train_labels
1446
                   # 3. Define the Objective Function for Optuna Hyperparameter Tuning
1447
                   def objective(trial):
                           """Defines the hyperparameter search space and evaluation metric for Optuna to optimize."""
1448
                           # Define the hyperparameter search space for the LightGBM model.
1449
                                  'objective': 'binary'
                                 'objective': 'binary',
'metric': 'auc',
'boosting_type': 'gbdt',
'n_estimators': trial.suggest_int('n_estimators', 2000, 10000),
'learning_rate': trial.suggest_float('learning_rate', 0.01, 0.05),
1450
1451
1452
                                 'num_leaves': trial.suggest_int('num_leaves', 20, 100),
'max_depth': trial.suggest_int('max_depth', 5, 10),
1453
                                  'seed': 42,
1454
                                 'n_jobs': -1,
'verbose': -1,
1455
                                  'colsample_bytree': trial.suggest_float('colsample_bytree', 0.5, 0.9),
1456
                                  'subsample': trial.suggest_float('subsample', 0.5, 0.9),
'reg_alpha': trial.suggest_float('reg_alpha', 0.01, 0.5),
1457
                                  'reg_lambda': trial.suggest_float('reg_lambda', 0.01, 0.5),
```

```
1458
1459
               N SPLITS = 5
               skf = StratifiedKFold(n splits=N SPLITS, shuffle=True, random state=42)
1460
               oof_auc_scores = []
1461
               # Perform cross-validation to get a robust estimate of the model's performance for a given set of
1462
                   hyperparameters
               for fold, (train_idx, val_idx) in enumerate(skf.split(X, y)):
    X_train, y_train = X.iloc[train_idx], y.iloc[train_idx]
1463
                   X_val, y_val = X.iloc[val_idx], y.iloc[val_idx]
1464
                   model = lgb.LGBMClassifier(**params)
1465
                   model.fit(X_train, y_train,
1466
                              eval_set=[(X_val, y_val)],
eval_metric='auc',
1467
                              callbacks=[lgb.early_stopping(100, verbose=False)])
1468
                   val_preds = model.predict_proba(X_val)[:, 1]
1469
                    auc = roc_auc_score(y_val, val_preds)
                   oof_auc_scores.append(auc)
1470
                # Return the mean AUC score, which Optuna will try to maximize.
1471
               return np.mean(oof_auc_scores)
1472
           # 4. Run the Optimization Study
1473
           print("Running hyperparameter optimization...")
           study = optuna.create_study(direction='maximize')
1474
           # The number of trials can be adjusted; 50 provides a reasonable balance between search time and
               performance.
1475
           study.optimize(objective, n_trials=50)
1476
           best_params = study.best_params
           print(f"\nBest trial parameters: {best_params}")
1477
           print(f"Best trial mean OOF AUC: {study.best_value}")
1478
           # 5. Train Final Model with Best Parameters and Create Submission
1479
            print("\nTraining final model with tuned parameters...")
           N_SPLITS = 5
1480
           skf = StratifiedKFold(n_splits=N_SPLITS, shuffle=True, random_state=42)
           test_preds = np.zeros(len(test_df))
oof_auc_scores = []
1481
1482
           # Combine the best parameters found by Optuna with fixed model settings.
           final_params = {
                'objective': 'binary', 'metric': 'auc', 'boosting_type': 'gbdt', 'seed': 42, 'n_jobs': -1, 'verbose': -1, **best_params
1484
1485
1486
           for fold, (train idx, val idx) in enumerate(skf.split(X, y)):
               print (f"===== Fold {fold+1} ====="
1487
               X_train, y_train = X.iloc[train_idx], y.iloc[train_idx]
               X_val, y_val = X.iloc[val_idx], y.iloc[val_idx]
1488
               model = lgb.LGBMClassifier(**final params)
1489
               1490
1491
                          callbacks=[lgb.early_stopping(100, verbose=False)])
1492
               val_preds = model.predict_proba(X_val)[:, 1]
               auc = roc_auc_score(y_val, val_preds)
1493
               oof_auc_scores.append(auc)
print(f"Fold {fold+1} AUC: {auc}")
1494
1495
                # Accumulate predictions on the test set, averaging over the folds to improve robustness.
               test_preds += model.predict_proba(test_df)[:, 1] / N_SPLITS
1496
1497
           # Print the final average validation score from the cross-validation.
           print(f"\nAverage OOF AUC with tuned parameters: {np.mean(oof_auc_scores)}")
1498
           # Create the submission DataFrame using the averaged test predictions.
1499
           submission_df = pd.DataFrame({'id': test_ids, 'Exited': test_preds})
           # Save the DataFrame to a CSV file in the 'final' directory for submission.submission_df.to_csv('final/submission.csv', index=False)
1501
1502
           print("\nSubmission file created successfully at 'final/submission.csv'.")
           print("First 5 rows of the submission file:")
1503
           print(submission df.head())
1504
```

Listing 15: Solution code and machine learning model.

# E.6 VISUALIZATION

1512

1513

1514

1515

1516

1517

1518 1519

1520 1521

```
f'''Draw a stacked bar chart that displays the percentage of restaurants of offering online ordering options versus those not offering online ordering options across different rating levels. Save this chart as result.jpg'. The title of the image should be "Percentage of Restaurants' Online Order Option by Rating", with the xlabel as "Rating", and the ylabel as "Percentage of Online Orders".'''
```

#### Listing 16: User's question.

```
1522
             import pandas as pd
             import sqlite3
1523
             import zipfile
1524
             import os
             import numpy as np
1525
             import matplotlib.pyplot as plt
1526
             def generate_chart():
1527
                 Loads data from the Zomato SQLite database, cleans the 'rate' column, calculates the percentage of restaurants offering online ordering for each rating,
1528
                  and generates a stacked bar chart of the results, saving it to the 'final' directory.
1529
1530
                 zip_path = 'data/zomato_rawdata.sqlite.zip'
db_name = 'zomato_rawdata.sqlite'
1531
                 extract_path = '.' # Extract to the current directory
output_dir = 'final'
1532
                 output_path = os.path.join(output_dir, 'result.jpg')
1533
                  # --- Step 1: Create output directory if it doesn't exist ---
1534
                 if not os.path.exists(output_dir):
1535
                      os.makedirs(output_dir)
1536
                    --- Step 2: Extract the database from the zip file ---
                 db_path = os.path.join(extract_path, db_name)
1537
                 with zipfile.ZipFile(zip_path, 'r') as zip_ref:

# Extract the database file, ensuring the correct member is chosen
                       for member in zip_ref.infolist():
1539
                           if db_name in member.filename:
                                 # To handle cases where the zip file has directory structures
                                member.filename = os.path.basename(member.filename)
                                zip ref.extract (member, path=extract path)
1541
1542
                  # --- Step 3: Load data from the SQLite database ---
1543
                  conn = None
1544
                       conn = sqlite3.connect(db_path)
                      # Select only the columns needed for this analysis query = "SELECT rate, online_order FROM Users"
1545
1546
                       df = pd.read_sql_query(query, conn)
                  finally:
1547
                      if conn:
1548
                           conn.close()
1549
                  # --- Step 4: Clean the 'rate' column -
                  # Make a copy to avoid SettingWithCopyWarning
1550
                 df = df.copv()
1551
                 # Replace special non-numeric strings 'NEW' and '-' with numpy.nan
df['rate'] = df['rate'].replace(['NEW', '-'], np.nan)
1552
1553
                 # Remove the '/5' suffix from the rating strings
df['rate'] = df['rate'].str.replace('/5', '').str.strip()
1554
                 # Convert the cleaned column to a numeric (float) type
df['rate'] = pd.to_numeric(df['rate'], errors='coerce')
1555
1556
                 # Drop rows where 'rate' or 'online_order' is NaN as they are not useful for this analysis
df.dropna(subset=['rate', 'online_order'], inplace=True)
1557
1558
                  # --- Step 5: Aggregate the data for plotting ---
# Create a cross-tabulation to get counts of 'Yes' and 'No' for each rating
1559
                 online_order_counts = pd.crosstab(df['rate'], df['online_order'])
1560
                   Calculate the percentage by dividing each row by its sum
1561
                  rate_online_percentage_df = online_order_counts.div(online_order_counts.sum(axis=1), axis=0) * 100
1562
                         Step 6: Generate and save the stacked bar chart -
1563
                 # Create the plot figure and axes
fig, ax = plt.subplots(figsize=(14, 8))
1564
                  # Generate the stacked bar chart
1565
                 rate_online_percentage_df.plot(kind='bar', stacked=True, ax=ax)
```

```
1566
                  # Set the title and labels as required
ax.set_title("Percentage of Restaurants' Online Order Option by Rating", fontsize=16)
ax.set_xlabel("Rating", fontsize=12)
1567
1568
                  ax.set_ylabel("Percentage of Online Orders", fontsize=12)
1569
                  # Improve aesthetics and readability
                  ax.tick_params(axis='x', rotation=90)
ax.legend(title='Online Order')
ax.grid(axis='y', linestyle='--', alpha=0.7)
1570
1571
1572
                 # Adjust layout to prevent labels from being cut off
plt.tight_layout()
1573
1574
                 # Save the resulting chart to the final directory
plt.savefig(output_path)
1575
                  # Close the plot to free up memory
1576
                  plt.close(fig)
1577
                  print(f"Chart has been generated and saved as '{output_path}'.")
1578
                   \# --- Step 7: Clean up the extracted file ---
1579
                  if os.path.exists(db_path):
                       os.remove(db_path)
1580
1581
             # Execute the function to generate the chart
1582
             generate_chart()
1583
```

Listing 17: Solution code.

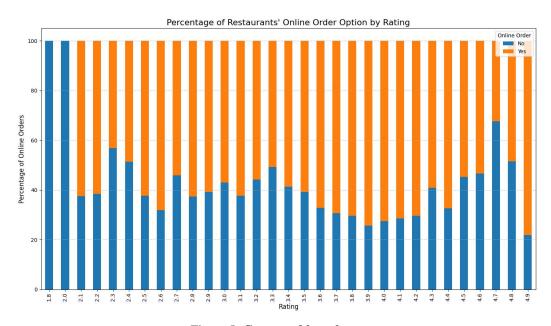


Figure 5: Generated bar chart.

# F FULL EXECUTION LOGS OF DS-STAR

In this section, we provide full logs of DS-STAR to better understand its behavior.

```
f'''In February 2023 what delta would Rafa_AI pay if the relative fee of \hookrightarrow the fee with ID=17 changed to 99?'''
```

Listing 18: User's question.

```
# Step 1
f'''First, I will filter the `payments.csv` data to select only the

→ transactions for the merchant 'Rafa_AI' that occurred in February

→ 2023. Since 2023 is not a leap year, February corresponds to the

→ days of the year from 32 to 59.'''
```

Listing 19: Round 0: Initial plan.

```
import pandas as pd

# Load the payments data
try:
    payments_df = pd.read_csv('data/payments.csv')
except FileNotFoundError:
    print("Error: 'data/payments.csv' not found. Please ensure the file is in the correct directory.")
    exit()

# Define the filter for the merchant 'Rafa_AI'
merchant_filter = payments_df['merchant'] == 'Rafa_AI'

# Define the filter for February 2023.
# 2023 is not a leap year. January has 31 days.
# February runs from day 32 (31+1) to day 59 (31+28).
date_filter = (payments_df['day_of_year'] >= 32) & (payments_df['day_of_year'] <= 59)

# Apply both filters to the DataFrame
rafa_ai_feb_transactions = payments_df[merchant_filter & date_filter]

# Display the resulting filtered DataFrame
print(rafa_ai_feb_transactions)</pre>
```

Listing 20: Round 0: Initial implementation.

```
f'''Insufficient'''
```

Listing 21: Round 0: Validation.

```
f'''Add Step'''
```

Listing 22: **Round 0: Routing.** Add Step 2 to the initial plan.

1689

```
1674
1675
       # Step 1
1676
       f'''First, I will filter the `payments.csv` data to select only the

→ transactions for the merchant 'Rafa_AI' that occurred in February

        → 2023. Since 2023 is not a leap year, February corresponds to the
1678
       \hookrightarrow days of the year from 32 to 59.'''
1679
1680
       # Step 2
1681
       f'''I will calculate the original total fee paid by Rafa_AI in February
1682
           2023. This involves loading the `merchant_data.json` and `fees.json`
           files, enriching the transaction data with all attributes required
          for fee matching (including merchant details, monthly volume,
1684
        \hookrightarrow monthly fraud rate, and intracountry status), then applying the
1685
       → corresponding fee rule to each transaction and summing the
1686
       \hookrightarrow calculated fees..'''
1687
```

#### Listing 23: Round 1: Plan.

```
1690
           import pandas as pd
1691
           import json
1692
           # --- 1. Load All Necessary Data ---
1693
               payments df = pd.read csv('data/payments.csv')
               merchant_data_df = pd.read_json('data/merchant_data.json')
fees_df = pd.read_json('data/fees.json')
1695
           acquirer_countries_df = pd.read_csv('data/acquirer_countries.csv')
except FileNotFoundError as e:
1696
               print(f"Error loading data: {e}. Please ensure all data files are in the 'data/' directory.")
1697
               exit()
             --- 2. Filter Transactions (Based on Previous Plan) ---
           # Define the filter for the merchant 'Rafa AI
1699
           merchant_filter = payments_df['merchant'] == 'Rafa_AI'
           # Define the filter for February 2023.
1701
           \# 2023 is not a leap year. January has 31 days. \# February runs from day 32 (31+1) to day 59 (31+28).
1702
           \label{linear} \verb|date_filter = (payments_df['day_of_year'] >= 32) & (payments_df['day_of_year'] <= 59) \\
1703
           # Apply both filters to the DataFrame and create a copy to avoid warnings
1704
           rafa_ai_feb_transactions = payments_df[merchant_filter & date_filter].copy()
1705
           # --- 3. Enrich Transaction Data with Merchant and Acquirer Information ---
1706
           # Get Rafa_AI's specific merchant data
1707
           rafa_ai_merchant_info = merchant_data_df[merchant_data_df['merchant'] == 'Rafa_AI'].iloc[0]
1708
           # Add merchant-specific attributes to the transaction DataFrame
           rafa_ai_feb_transactions['account_type'] = rafa_ai_merchant_info['account_type']
rafa_ai_feb_transactions['capture_delay'] = str(rafa_ai_merchant_info['capture_delay']) # Ensure consistent
1709
1710
           rafa_ai_feb_transactions['merchant_category_code'] = rafa_ai_merchant_info['merchant_category_code']
1711
           # FIX: The original code failed on the next line.
1712
           # The error "ValueError: Length of values (1) does not match length of index (2126)"
             suggests that rafa_ai_merchant_info['acquirer'] is a list-like object with one element
1713
             (e.g., ['gringotts']) instead of a scalar string. Assigning a list of length 1 to a
1714
           # column of a different length causes this error.
           # The fix is to extract the scalar string from the list before assigning it.
1715
           acquirer_name = rafa_ai_merchant_info['acquirer']
           if isinstance(acquirer_name, list):
    # If it's a list, take the first element. If empty, it becomes None.
1716
               acquirer_name = acquirer_name[0] if acquirer_name else None
1717
1718
           # Now assign the scalar acquirer name to the new column.
           rafa_ai_feb_transactions['acquirer'] = acquirer_name
1719
1720
           # Create a mapping from acquirer name to country code
1721
           acquirer_country_map = acquirer_countries_df.set_index('acquirer')['country_code']
1722
           # Add acquirer country to the transaction data
           rafa_ai_feb_transactions['acquirer_country']
1723
               rafa_ai_feb_transactions['acquirer'].map(acquirer_country_map)
1724
            Determine if each transaction is intracountry
1725
           1726
1727
           # --- 4. Calculate Monthly Aggregates for Fee Matching ---
           # Calculate total transaction volume for Rafa_AI in February
```

```
1728
            monthly_volume = rafa_ai_feb_transactions['eur_amount'].sum()
1729
            # Calculate fraudulent transaction volume for the same period
1730
            fraudulent_volume
1731

→ rafa_ai_feb_transactions[rafa_ai_feb_transactions['has_fraudulent_dispute']]['eur_amount'].sum()

1732
            # Calculate the monthly fraud rate as a percentage
            # Avoid division by zero if there's no volume
monthly_fraud_rate = (fraudulent_volume / monthly_volume) if monthly_volume > 0 else 0
1733
            monthly_fraud_percentage = monthly_fraud_rate * 100
1734
1735
            # Add these aggregates to each transaction row for easy lookup during fee matching
            rafa_ai_feb_transactions['monthly_volume'] = monthly_volume
1736
            rafa_ai_feb_transactions['monthly_fraud_percentage'] = monthly_fraud_percentage
1737
1738
               --- 5. Implement Fee Matching Logic ---
1739
            def parse_numeric_value(s):
                    'Helper to parse strings like '100k' or '1m' into numbers."""
1740
                s = s.strip().lower()
if 'k' in s:
1741
                     return float(s.replace('k', '')) * 1_000
1742
                if 'm' in s:
                     return float(s.replace('m', '')) * 1_000_000
1743
                return float(s)
1744
            def check_volume_rule(rule_str, value):
                     Checks if a value fits a volume rule string (e.g., '100k-1m')."""
1745
                if pd.isnull(rule_str): return True
1746
                rule_str = rule_str.replace(' ', '
if '-' in rule str.
                        in rule_str:
1747
                     low_str, high_str = rule_str.split('-')
                return parse_numeric_value(low_str) <= value < parse_numeric_value(high_str)
if '>' in rule_str:
1748
1749
                     return value >= parse_numeric_value(rule_str.replace('>', ''))
                if '<' in rule_str:</pre>
1750
                     return value < parse_numeric_value(rule_str.replace('<', ''))</pre>
                return False
1751
1752
            def check_fraud_rule(rule_str, value):
                """Checks if a value fits a fraud rule string (e.g., '7.7%-8.3%').""" if pd.isnull(rule_str): return True
1753
                 rule_str = rule_str.replace('%', '')
1754
                       in rule str:
                     low_str, high_str = rule_str.split('-')
1755
                     return float(low_str) <= value < float(high_str)
                if '>' in rule_str:
1756
                     return value >= float(rule_str.replace('>', ''))
1757
                if '<' in rule str:
                     return value < float(rule_str.replace('<', ''))
1758
                return False
1759
            def find_matching_fee_rule(transaction, all_rules_df):
1760
                  ""Iterates through fee rules to find the first one that matches the transaction."""
                for _, rule in all_rules_df.iterrows():
1761
                     # A rule with an empty list for a category applies to all values in that category if rule['card_scheme']!= transaction['card_scheme']: continue if rule['account_type'] and transaction['account_type'] not in rule['account_type']: continue
1762
1763
                     if rule['merchant_category_code'] and transaction['merchant_category_code'] not in

→ rule['merchant_category_code']: continue
if rule['aci'] and transaction['aci'] not in rule['aci']: continue

1764
1765
                     # A rule with a null value for a category applies to all values
if pd.notnull(rule['is_credit']) and rule['is_credit'] != transaction['is_credit']: continue
1766
                     if pd.notnull(rule['intracountry']) and rule['intracountry'] != transaction['intracountry']:
1767
                         continue
                     if pd.notnull(rule['capture_delay']) and rule['capture_delay'] != transaction['capture_delay']:
1768
1769
                     # Check complex range-based rules
1770
                     if not check_volume_rule(rule['monthly_volume'], transaction['monthly_volume']): continue
if not check_fraud_rule(rule['monthly_fraud_level'], transaction['monthly_fraud_percentage']):
1771
1772
                     # If all checks pass, this is the correct rule
                     return rule['fixed_amount'], rule['rate']
1773
1774
                 # Return default/zero fees if no rule is found (should not happen with complete data)
                return 0, 0
1775
                -- 6. Calculate Fees and Final Sum -
1776
1777
            # Apply the matching function to each transaction to get its specific fee parameters
            fee_params = rafa_ai_feb_transactions.apply(
1778
                lambda row: find_matching_fee_rule(row, fees_df),
                axis=1
1779
1780
            # Assign the found fee parameters back to the DataFrame
1781
            rafa_ai_feb_transactions[['fixed_amount', 'rate']] = pd.DataFrame(
                fee params.tolist(),
```

17951796

179717981799

1800 1801

1802

1822

```
1782
                index=rafa_ai_feb_transactions.index
1783
1784
             Calculate the fee for each transaction using the formula from the manual
            # fee = fixed_amount + rate * transaction_value / 10000
rafa_ai_feb_transactions['fee'] = rafa_ai_feb_transactions['fixed_amount'] + ''
1785
1786
                                                 (rafa_ai_feb_transactions['rate'] * rafa_ai_feb_transactions['eur_amount'] /
                                                \hookrightarrow 10000)
1787
            \ensuremath{\sharp} Calculate the total fee by summing up the fees for all transactions
1788
            total_original_fee = rafa_ai_feb_transactions['fee'].sum()
1789
            # Display the final result
1790
            print(f"The original total fee paid by Rafa_AI in February 2023 was: {total_original_fee:.2f} EUR")
1791
```

Listing 24: Round 1: Implementation.

```
f'''Insufficient'''
```

Listing 25: **Round 1: Validation.** 

```
f'''Add Step'''
```

Listing 26: **Round 1: Routing.** Add Step 3 to the current plan.

```
1803
         Step 1
1804
       f'''First, I will filter the `payments.csv` data to select only the
1805
           transactions for the merchant 'Rafa_AI' that occurred in February
        \hookrightarrow 2023. Since 2023 is not a leap year, February corresponds to the
1807
        \hookrightarrow days of the year from 32 to 59.'''
       # Step 2
1809
       f'''I will calculate the original total fee paid by Rafa_AI in February
1810
        \leftrightarrow 2023. This involves loading the `merchant_data.json` and `fees.json`
1811
           files, enriching the transaction data with all attributes required
1812
           for fee matching (including merchant details, monthly volume,
           monthly fraud rate, and intracountry status), then applying the
1813
           corresponding fee rule to each transaction and summing the
        \hookrightarrow
1814
           calculated fees..'''
1815
1816
       # Step 3
1817
       f'''I will calculate the new total fee for Rafa_AI's February 2023
1818
           transactions using the modified fee structure where the fee rule
           with ID=17 has its `rate` changed to 99. Then, I will subtract the
1819
           original total fee (778.52 EUR) from this new total fee to find the
        \hookrightarrow
1820
           requested delta.'''
1821
```

Listing 27: Round 2: Plan.

```
import pandas as pd
            import ison
1826
            # --- 1. Load All Necessary Data ---
1827
                 payments_df = pd.read_csv('data/payments.csv')
1828
                 merchant_data_df = pd.read_json('data/merchant_data.json')
fees_df = pd.read_json('data/fees.json')
1829
            acquirer_countries.df = pd.read_csv('data/acquirer_countries.csv')
except FileNotFoundError as e:
1830
                 print(f"Error loading data: {e}. Please ensure all data files are in the 'data/' directory.")
1831
                 exit()
1832
                   2. Filter Transactions (Based on Previous Plan) ---
            # Define the filter for the merchant 'Rafa_AI'
1833
            merchant_filter = payments_df['merchant'] == 'Rafa_AI'
1834
            # Define the filter for February 2023.
1835
            # 2023 is not a leap year. January has 31 days.
# February runs from day 32 (31+1) to day 59 (31+28).
```

```
1836
            date_filter = (payments_df['day_of_year'] >= 32) & (payments_df['day_of_year'] <= 59)</pre>
1837
            # Apply both filters to the DataFrame and create a copy to avoid warnings
1838
            rafa_ai_feb_transactions = payments_df[merchant_filter & date_filter].copy()
1839
            # --- 3. Enrich Transaction Data with Merchant and Acquirer Information --
1840
            # Get Rafa_AI's specific merchant data
1841
            rafa_ai_merchant_info = merchant_data_df[merchant_data_df['merchant'] == 'Rafa_AI'].iloc[0]
1842
            # Add merchant-specific attributes to the transaction DataFrame
rafa_ai_feb_transactions['account_type'] = rafa_ai_merchant_info['account_type']
1843
            rafa_ai_feb_transactions['capture_delay'] = str(rafa_ai_merchant_info['capture_delay']) # Ensure consistent
1844

→ type

            rafa_ai_feb_transactions['merchant_category_code'] = rafa_ai_merchant_info['merchant_category_code']
1845
            # FIX: The original code failed on the next line.
1846
            # The error "ValueError: Length of values (1) does not match length of index (2126)"
1847
            # suggests that rafa_ai_merchant_info['acquirer'] is a list-like object with one element
# (e.g., ['gringotts']) instead of a scalar string. Assigning a list of length 1 to a
1848
              column of a different length causes this error.
            # The fix is to extract the scalar string from the list before assigning it.
acquirer_name = rafa_ai_merchant_info['acquirer']
1849
1850
            if isinstance(acquirer_name, list):
                # If it's a list, take the first element. If empty, it becomes None.
1851
                acquirer_name = acquirer_name[0] if acquirer_name else None
1852
            # Now assign the scalar acquirer name to the new column.
1853
            rafa_ai_feb_transactions['acquirer'] = acquirer_name
1854
            # Create a mapping from acquirer name to country code
1855
            acquirer_country_map = acquirer_countries_df.set_index('acquirer')['country_code']
1856
             Add acquirer country to the transaction data
1857
            rafa_ai_feb_transactions['acquirer_country'] =
    rafa_ai_feb_transactions['acquirer'].map(acquirer_country_map)
1858
            # Determine if each transaction is intracountry
rafa_ai_feb_transactions['intracountry'] = rafa_ai_feb_transactions['issuing_country'] ==
1859
             → rafa_ai_feb_transactions['acquirer_country']
1861
            # --- 4. Calculate Monthly Aggregates for Fee Matching -
1862
            # Calculate total transaction volume for Rafa AI in February
            monthly_volume = rafa_ai_feb_transactions['eur_amount'].sum()
1863
1864
            # Calculate fraudulent transaction volume for the same period
1865

→ rafa_ai_feb_transactions[rafa_ai_feb_transactions['has_fraudulent_dispute']]['eur_amount'].sum()

1866
            # Calculate the monthly fraud rate as a percentage
            # Avoid division by zero if there's no volume
monthly_fraud_rate = (fraudulent_volume / monthly_volume) if monthly_volume > 0 else 0
1867
1868
            monthly_fraud_percentage = monthly_fraud_rate * 100
1869
             Add these aggregates to each transaction row for easy lookup during fee matching
            rafa_ai_feb_transactions['monthly_volume'] = monthly_volume
rafa_ai_feb_transactions['monthly_fraud_percentage'] = monthly_fraud_percentage
1870
1871
1872
            # --- 5. Implement Fee Matching Logic ---
1873
            def parse_numeric_value(s):
                  ""Helper to parse strings like '100k' or '1m' into numbers."""
1874
                 s = s.strip().lower()
1875
                if 'k' in s:
                     return float(s.replace('k', '')) * 1_000
1876
                if 'm' in s:
                    return float(s.replace('m', '')) * 1_000_000
1877
                return float(s)
1878
            def check volume rule (rule str. value):
1879
                 """Checks if a value fits a volume rule string (e.g., '100k-1m')."""
                if pd.isnull(rule_str): return True
rule_str = rule_str.replace(' ', ''
1880
                if '-' in rule_str:
1881
                     low_str, high_str = rule_str.split('-')
1882
                     return parse_numeric_value(low_str) <= value < parse_numeric_value(high_str)</pre>
                if '>' in rule_str:
1883
                return value >= parse_numeric_value(rule_str.replace('>', ''))
if '<' in rule_str:</pre>
1884
                    return value < parse_numeric_value(rule_str.replace('<', ''))</pre>
1885
                return False
1886
            def check_fraud_rule(rule_str, value):
    """Checks if a value fits a fraud rule string (e.g., '7.7%-8.3%')."""
1887
                 if pd.isnull(rule_str): return True
1888
                rule_str = rule_str.replace('%',
if '-' in rule_str:
1889
                     low_str, high_str = rule_str.split('-')
                     return float(low_str) <= value < float(high_str)</pre>
```

```
1890
                 if '>' in rule_str:
1891
                     return value >= float(rule_str.replace('>', ''))
                 if '<' in rule_str:</pre>
1892
                     return value < float(rule_str.replace('<', ''))</pre>
1893
                 return False
1894
            def find_matching_fee_rule(transaction, all_rules_df):
                 """Iterates through fee rules to find the first one that matches the transaction."""
for _, rule in all_rules_df.iterrows():
1895
                     # A rule with an empty list for a category applies to all values in that category if rule['card_scheme'] != transaction['card_scheme']: continue if rule['account_type'] and transaction['account_type'] not in rule['account_type']: continue
1896
1897
                     if rule['merchant_category_code'] and transaction['merchant_category_code'] not in 

→ rule['merchant_category_code']: continue 
if rule['aci'] and transaction['aci'] not in rule['aci']: continue
1898
                     # A rule with a null value for a category applies to all values
if pd.notnull(rule['is_credit']) and rule['is_credit'] != transaction['is_credit']: continue
if pd.notnull(rule['intracountry']) and rule['intracountry'] != transaction['intracountry']:
1900
1901
                          continue
1902
                      if pd.notnull(rule['capture_delay']) and rule['capture_delay'] != transaction['capture_delay']:
1903
1904
                     if not check_volume_rule(rule['monthly_volume'], transaction['monthly_volume']): continue
if not check_fraud_rule(rule['monthly_fraud_level'], transaction['monthly_fraud_percentage']):
1905
1906
1907
                      # If all checks pass, this is the correct rule
                     return rule['fixed_amount'], rule['rate']
1908
                 # Return default/zero fees if no rule is found (should not happen with complete data)
1909
                 return 0. 0
1910
            # --- 6. Calculate Fees and Final Sum ---
1911
            # Apply the matching function to each transaction to get its specific fee parameters
1912
            fee_params = rafa_ai_feb_transactions.apply(
                 lambda row: find_matching_fee_rule(row, fees_df),
1913
1914
1915
            # Assign the found fee parameters back to the DataFrame
            rafa_ai_feb_transactions[['fixed_amount', 'rate']] = pd.DataFrame(
1916
                 fee params.tolist(),
                 index=rafa_ai_feb_transactions.index
1917
1918
            \sharp Calculate the fee for each transaction using the formula from the manual
1919
            # fee = fixed_amount + rate * transaction_value / 10000
            rafa_ai_feb_transactions['fee'] = rafa_ai_feb_transactions['fixed_amount'] + '
1920
                                                  (rafa_ai_feb_transactions['rate'] * rafa_ai_feb_transactions['eur_amount'] /
1921
                                                  \hookrightarrow 10000)
1922
            # Calculate the total fee by summing up the fees for all transactions
total_original_fee = rafa_ai_feb_transactions['fee'].sum()
1923
            # Display the final result from the previous step
print(f"The original total fee paid by Rafa_AI in February 2023 was: {total_original_fee:.2f} EUR")
1924
1925
1926
            # --- 7. Calculate New Fee with Modified Rule ---
1927
            # Create a copy of the fees dataframe to avoid altering the original
            fees_df_modified = fees_df.copy()
1928
1929
            # Implement the plan: change the rate for the rule with ID=17 to 99
            fees_df_modified.loc[fees_df_modified['ID'] == 17, 'rate'] = 99
1930
            # Apply the matching function again, this time with the modified fee rules
1931
            new_fee_params = rafa_ai_feb_transactions.apply(
1932
                 lambda row: find_matching_fee_rule(row, fees_df_modified),
                 axis=1
1933
1934
            # Assign the new fee parameters to new columns in the DataFrame
            rafa_ai_feb_transactions[['new_fixed_amount', 'new_rate']] = pd.DataFrame(
1935
                 new fee params.tolist(),
1936
                 index=rafa_ai_feb_transactions.index
1937
             # Calculate the new fee for each transaction using the new rate
1938
            1939
                                                         \hookrightarrow rafa_ai_feb_transactions['eur_amount'] / 10000)
1940
            # Calculate the new total fee by summing up the individual new fees
1941
            new total fee = rafa ai feb transactions['new fee'].sum()
1942
            # Calculate the difference (delta) between the new total fee and the original total fee
1943
            fee_delta = new_total_fee - total_original_fee
```

1945

1946

1948

1950

1951

1952 1953 1954

1955 1956

1957 1958 1959

1960

1961

1962

1963

1965

1966

1967

1969

1971

1973

1974

1975

1976

1977 1978

```
8. Display the Final Delta -
print(f"The new total fee with the modified rule would be: {new_total_fee:.2f} EUR")
print(f"The requested delta (New Fee - Original Fee) is: {fee_delta:.2f} EUR")
```

Listing 28: **Round 2: Implementation.** 

```
f'''Insufficient'''
```

Listing 29: Round 2: Validation.

```
f'''Step 3 is wrong'''
```

Listing 30: Round 2: Routing. Remove Step 3 from the current plan and add Step 3 again.

```
# Step 1
       f'''First, I will filter the `payments.csv` data to select only the
          transactions for the merchant 'Rafa_AI' that occurred in February
       \hookrightarrow 2023. Since 2023 is not a leap year, February corresponds to the
          days of the year from 32 to 59.'''
       # Step 2
       f'''I will calculate the original total fee paid by Rafa_AI in February
          2023. This involves loading the `merchant_data.json` and `fees.json`
           files, enriching the transaction data with all attributes required
           for fee matching (including merchant details, monthly volume,
          monthly fraud rate, and intracountry status), then applying the
          corresponding fee rule to each transaction and summing the
       \hookrightarrow
          calculated fees..'''
       # Step 3
1972
       f'''I will filter the transactions for 'Rafa_AI' in February 2023
          (day_of_year from 32 to 59) that meet all the conditions of fee rule
           ID=17: `card_scheme` must be 'SwiftCharge', `is_credit` must be
           True, and `aci` must be 'A'. Then, for these matching transactions,
           I will calculate the total fee difference using the formula: (99 -
           60) * SUM(eur_amount) / 10000.'''
```

Listing 31: Round 3: Plan.

```
1981
              import pandas as pd
1982
               --- 1. Load All Necessary Data -
1983
             # This part is based on the loading logic from the base code.
1984
                  payments_df = pd.read_csv('data/payments.csv')
              merchant_data_df = pd.read_json('data/merchant_data.json')
fees_df = pd.read_json('data/fees.json')
except FileNotFoundError as e:
                  print(f"Error loading data: {e}. Please ensure all data files are in the 'data/' directory.")
1987
                  exit()
1988
              # --- Implementation of the Current Plan ---
# I will filter the transactions for 'Rafa_AI' in February 2023 (day_of_year from 32 to 59)
1989
             # that meet all the conditions of fee rule ID=17: `card_scheme` must be 'SwiftCharge',
# `is_credit` must be True, and `aci` must be 'A'. Then, for these matching transactions,
# I will calculate the total fee difference using the formula:
1991
             # (99 - 60) * SUM(eur_amount) / 10000.
1992
             print("--- Calculating fee difference for 'Rafa_AI' in Feb 2023 under specific conditions ---")
1993
             # Step 1: Define the conditions for filtering based on the plan.
1994
               February 2023 corresponds to day_of_year 32 to 59 (2023 is not a leap year).
1995
             feb_start_day = 32
feb_end_day = 59
1996
             merchant_name = 'Rafa_AI'
             required_card_scheme = 'SwiftCharge'
1997
              required_is_credit = True
             required_aci = 'A'
```

202020212022

2023

202420252026

20272028

202920302031

2032

2033

2034

2035

2037

2038

2039

2040

2042

2043

2044

2045

2046

2047 2048

```
1998
1999
             \# Step 2: Filter the payments DataFrame to find transactions matching all conditions.
            matching_transactions = payments_df[
2000
                  (payments_df['merchant'] == merchant_name) & (payments_df['day_of_year'] >= feb_start_day) & (payments_df['day_of_year'] <= feb_end_day) &
2001
2002
                  (payments_df['card_scheme'] == required_card_scheme) &
                  (payments_df['is_credit'] == required_is_credit) &
(payments_df['aci'] == required_aci)
2003
2004
2005
            print(f"\n[INFO] Found {len(matching_transactions)} transactions for '{merchant_name}' in February 2023
            \hookrightarrow matching all specified conditions.")
2006
             # Step 3: Calculate the total fee difference for these matching transactions.
            if not matching_transactions.empty:
    # Sum the 'eur_amount' for all matching transactions.
    total_matching_volume = matching_transactions['eur_amount'].sum()
2008
2009
                 print(f"[INFO] Total volume of these transactions: {total_matching_volume:.2f} EUR")
2010
                  # Apply the fee difference formula provided in the plan
                 # Formula: (new_rate - old_rate) * SUM(eur_amount) / 10000
rate_difference = 99 - 60
2011
2012
                 total_fee_difference = (rate_difference * total_matching_volume) / 10000
2013
                     -- Final Conclusion ---
                                               -----CONCLUSION-----
2014
                 print (f"The total fee difference for the matching transactions is: {total_fee_difference:.4f} EUR.")
2015
2016
                 # --- Final Conclusion ---
                                                   ---CONCLUSION---
2017
                 print("No transactions were found that match all the specified criteria.")
                 print("Therefore, the total fee difference is 0.00 EUR.")
2018
```

Listing 32: Round 3: Implementation.

```
f'''Insufficient'''
```

Listing 33: **Round 3: Validation.** 

```
f'''Step 3 is wrong'''
```

Listing 34: Round 3: Routing. Remove Step 3 from the current plan and add Step 3 again.

```
# Step 1
f'''First, I will filter the `payments.csv` data to select only the
→ transactions for the merchant 'Rafa_AI' that occurred in February
→ 2023. Since 2023 is not a leap year, February corresponds to the
\rightarrow days of the year from 32 to 59.'''
# Step 2
f'''I will calculate the original total fee paid by Rafa_AI in February
   2023. This involves loading the `merchant_data.json` and `fees.json`
   files, enriching the transaction data with all attributes required
\hookrightarrow for fee matching (including merchant details, monthly volume,
→ monthly fraud rate, and intracountry status), then applying the
→ corresponding fee rule to each transaction and summing the
→ calculated fees..'''
# Step 3
f'''Calculate the fee delta by subtracting the original `rate` of fee ID
   17 from the new `rate` of 99, multiplying the result by the total
   volume of the affected transactions (686.48 EUR), and then dividing
```

Listing 35: Round 4: Plan.

```
import pandas as pd import json
```

2086208720882089

2090 2091

209220932094

209520962097

```
2052
2053
            # --- 1. Load All Necessary Data ---
            # This part is based on the loading logic from the base code.
2054
                payments_df = pd.read_csv('data/payments.csv')
merchant_data_df = pd.read_json('data/merchant_data.json')
fees_df = pd.read_json('data/fees.json')
2055
2056
            except FileNotFoundError as e:
2057
                print(f"Error loading data: {e}. Please ensure all data files are in the 'data/' directory.")
                exit()
2058
              --- Implementation of the Current Plan -
2059
            \# Calculate the fee delta by subtracting the original `rate` of fee ID 17 from the new `rate` of 99,
2060
            \# multiplying the result by the total volume of the affected transactions (686.48 EUR), \# and then dividing by 10000.
2061
            print("--- Calculating fee delta for affected 'Rafa AI' transactions ---")
2062
2063
            \ensuremath{\mathtt{\#}} Step 1: Define the constants given in the plan.
            new_rate = 99
2064
            total_affected_volume = 686.48 # EUR
            fee\_divisor = 10000
2065
            # Step 2: Find the original rate for fee ID 17 from the fees_df DataFrame.
2066
            fee\_rule\_17 = fees\_df[fees\_df['ID'] == 17]
2067
            if fee_rule_17.empty:
2068
                print("\n[ERROR] Fee rule with ID 17 was not found in fees.json.")
                original_rate = None # Set to None to prevent further calculation
2069
2070
                # Extract the 'rate' value from the filtered DataFrame.
                    iloc[0] selects the first row of the result.
2071
                ".inot(o) selected first in the restate.
original_rate = fee_rule_17['rate'].iloc(0)
print(f"\n[INFO] Original rate for fee ID 17 is: {original_rate}")
print(f"[INFO] New rate to be applied is: {new_rate}")
2072
2073
                print(f"[INFO] Total volume of affected transactions is: {total_affected_volume:.2f} EUR")
2074
            # Step 3: Calculate the fee delta using the specified formula.
            if original_rate is not None:
    # Calculate the difference in rates.
2075
                rate_difference = new_rate - original_rate
2076
                # Apply the formula: (rate_difference * total_volume) / 10000
2077
                fee_delta = (rate_difference * total_affected_volume) / fee_divisor
2078
                # --- Final Conclusion ---
2079
                print("\n--
                                          -----")
                print(f"The calculated fee delta, by changing the rate from {original_rate} to {new_rate} for the
2080

→ affected volume, is: {fee_delta:.4f} EUR.")
2081
2082
                # --- Final Conclusion ---
                                         -----CONCLUSION-----
2083
                print("Could not perform calculation because the original fee rate for ID 17 was not found.")
2084
```

Listing 36: Round 4: Implementation.

```
f'''Insufficient'''
```

Listing 37: Round 4: Validation.

```
f'''Add Step'''
```

Listing 38: **Round 4: Routing.** Add Step 4 to the current plan.

```
2100
2101

# Step 1
f'''First, I will filter the `payments.csv` data to select only the

→ transactions for the merchant 'Rafa_AI' that occurred in February

→ 2023. Since 2023 is not a leap year, February corresponds to the

→ days of the year from 32 to 59.'''

# Step 2
```

```
2106
       f'''I will calculate the original total fee paid by Rafa_AI in February
2107
        \leftrightarrow 2023. This involves loading the 'merchant_data.json' and 'fees.json'
2108
        \,\hookrightarrow\, files, enriching the transaction data with all attributes required
        \hookrightarrow for fee matching (including merchant details, monthly volume,
2109
        → monthly fraud rate, and intracountry status), then applying the
2110
        → corresponding fee rule to each transaction and summing the
2111

    calculated fees..'''

2112
2113
        # Step 3
2114
       f'''Calculate the fee delta by subtracting the original `rate` of fee ID
        → 17 from the new `rate` of 99, multiplying the result by the total
2115
        \hookrightarrow volume of the affected transactions (686.48 EUR), and then dividing
2116
        2117
2118
       # Step 4
2119
       f'''Calculate the total original fee for Rafa_AI's transactions in
           February 2023 that are affected by fee ID=17, and then calculate the
2120
           new total fee for the same transactions using the modified rate of
2121
           99. The delta is the difference between these two amounts.'''
2122
```

Listing 39: Round 5: Plan.

```
2125
2126
             import pandas as pd
import json
2127
             \# --- 1. Load All Necessary Data --- \# This part is based on the loading logic from the base code.
2128
2129
                  payments df = pd.read csv('data/payments.csv')
2130
                  merchant_data_df = pd.read_json('data/merchant_data.json')
                  fees_df = pd.read_json('data/fees.json')
2131
             except FileNotFoundError as e:
                  print(f"Error loading data: {e}. Please ensure all data files are in the 'data/' directory.")
2132
                  exit()
2133
             # --- Implementation of the Current Plan --
2134
             \# Calculate the total original fee for Rafa_AI's transactions in February 2023
             # that are affected by fee ID=17, and then calculate the new total fee for the same # transactions using the modified rate of 99. The delta is the difference.
2135
2136
             print("--- Identifying transactions for 'Rafa AI' in February 2023 affected by fee rule ID 17 ---")
2137
             # Step 1: Filter payments for 'Rafa_AI' in February 2023.
# February corresponds to days of the year from 32 to 59 (since 2023 is not a leap year).
2138
             feb_transactions_df = payments_df[
   (payments_df['merchant'] == 'Rafa_AI') &
        (payments_df['day_of_year'] >= 32) &
2139
2140
                  (payments_df['day_of_year'] <= 59)</pre>
2141
             ].copy()
2142
             # Step 2: Enrich transaction data with merchant-specific information. # Get the specific data for 'Rafa_AI' from the merchant data file.
2143
             rafa_ai_data = merchant_data_df[merchant_data_df['merchant'] == 'Rafa_AI']
2144
             # Merge the February transactions with the merchant's data.
2145
             if not rafa_ai_data.empty and not feb_transactions_df.empty:
2146
                  # Use a left merge to add merchant details (like mcc, account_type) to each transaction.
full_transaction_data = feb_transactions_df.merge(rafa_ai_data, on='merchant', how='left')
2147
                  print("\n[ERROR] No transactions found for 'Rafa AI' in February or 'Rafa AI' data not found.")
2148
2149
             # Step 3: Retrieve the specific conditions for the fee rule with ID=17.
2150
             fee_rule_17 = fees_df[fees_df['ID'] == 17]
2151
             if fee_rule_17.empty:
                  print("\n[ERROR] Fee rule with ID 17 was not found in fees.json.")
2152
                  exit()
2153
             else:
                  # Use .iloc[0] to get the rule as a Series
fee_rule_17 = fee_rule_17.iloc[0]
2154
2155
             # Step 4: Filter the transactions that match all conditions of the fee rule.
# The manual states that a null or empty list in a rule field means it applies to all values.
2156
             mask = pd.Series(True, index=full transaction data.index)
2157
             # Match card_scheme (string)
2158
                  mask &= (full_transaction_data['card_scheme'] == fee_rule_17['card_scheme'])
2159
             # Match account_type (list)
```

219621972198

2199 2200

```
2160
           if fee_rule_17['account_type']: # An empty list evaluates to False, correctly skipping the filter
   mask &= full_transaction_data['account_type'].isin(fee_rule_17['account_type'])
2161
2162
            # Match merchant_category_code (list)
2163
            if fee_rule_17['merchant_category_code']:
                mask &= full transaction data['merchant category code'].isin(fee rule 17['merchant category code'])
2164
            # Match aci (list)
2165
            if fee_rule_17['aci']:
                mask &= full_transaction_data['aci'].isin(fee_rule_17['aci'])
2166
            # Match is_credit (boolean, but can be null)
if fee_rule_17['is_credit'] is not None:
2167
2168
                mask &= (full_transaction_data['is_credit'] == fee_rule_17['is_credit'])
2169
            # Apply the combined mask to get the affected transactions
affected_transactions_df = full_transaction_data[mask]
2170
2171
            \sharp Step 5: Calculate the original and new fees for the affected transactions.
            if affected_transactions_df.empty:
2172
                print("\n[INFO] No transactions for 'Rafa_AI' in February match the criteria for fee rule ID 17.")
                 total_original_fee = 0
2173
                 total_new_fee = 0
2174
                 fee delta = 0
            else:
2175
                total_affected_volume = affected_transactions_df['eur_amount'].sum()
                num_affected_transactions = len(affected_transactions_df)
2176
                print(f"\n[INFO] Found {num_affected_transactions} transactions matching the fee rule.")
2177
                print(f"[INFO] Total volume of these transactions: {total_affected_volume:.2f} EUR")
2178
                 # Get fee parameters from the rule
2179
                original_rate = fee_rule_17['rate']
fixed_amount = fee_rule_17['fixed_amount']
2180
                 new_rate = 99
2181
                fee_divisor = 10000
2182
                 \# Calculate total original fee using the formula: fee = fixed_amount + rate * value / 10000
                # Vectorized calculation for efficiency
total_fixed_component = fixed_amount * num_affected_transactions
2183
2184
                 total_original_rate_component = (original_rate * total_affected_volume) / fee_divisor
                total_original_fee = total_fixed_component + total_original_rate_component
2185
                 # Calculate total new fee with the modified rate
2186
                total_new_rate_component = (new_rate * total_affected_volume) / fee_divisor total_new_fee = total_fixed_component + total_new_rate_component
2187
2188
                 # The delta is the difference between the new and original total fees
                fee_delta = total_new_fee - total_original_fee
2189
2190
            # --- Final Conclusion ---
                                     -----CONCLUSION-----")
2191
            print("\n---
            print(f"Total original fee for affected transactions: {total_original_fee:.4f} EUR")
2192
            print(f"Total new fee with rate={new_rate}: {total_new_fee:.4f} EUR")
print(f"The calculated fee delta is: {fee_delta:.4f} EUR.")
2193
            print ("This represents the total increase in fees for Rafa_AI's affected February transactions due to the
            → rate change.")
2194
```

Listing 40: **Round 5: Implementation.** 

```
f'''Sufficient'''
```

Listing 41: **Round 4: Validation.** 

```
2202
               import pandas as pd
2203
              import ison
2204
               \# --- 1. Load All Necessary Data --- \# This part is based on the loading logic from the base code.
2205
              payments_df = pd.read_csv('data/payments.csv')
2206
              merchant_data_df = pd.read_json('data/merchant_data.json')
fees_df = pd.read_json('data/fees.json')
2207
                   -- Implementation of the Calculation -
               # Calculate the total original fee for Rafa_AI's transactions in February 2023
              # that are affected by fee ID=17, and then calculate the new total fee for the same
# transactions using the modified rate of 99. The delta is the difference.
2210
2211
              # Step 1: Filter payments for 'Rafa_AI' in February 2023.
# February corresponds to days of the year from 32 to 59 (since 2023 is not a leap year).
2212
               feb_transactions_df = payments_df[
                    (payments_df['merchant'] == 'Rafa_AI') &
(payments_df['day_of_year'] >= 32) &
2213
                     (payments_df['day_of_year'] <= 59)
```

```
2214
           ].copy()
2215
           # Step 2: Enrich transaction data with merchant-specific information.
2216
            # Get the specific data for 'Rafa_AI' from the merchant data file.
           rafa_ai_data = merchant_data_df[merchant_data_df['merchant'] == 'Rafa_AI']
2217
2218
           # Merge the February transactions with the merchant's data.
           full_transaction_data = pd.DataFrame()
2219
           if not rafa_ai_data.empty and not feb_transactions_df.empty:
    # Use a left merge to add merchant details (like mcc, account_type) to each transaction.
2220
               full_transaction_data = feb_transactions_df.merge(rafa_ai_data, on='merchant', how='left')
           # Step 3: Retrieve the specific conditions for the fee rule with ID=17.
2222
           fee_delta = 0.0
           if not full_transaction_data.empty:
2223
               fee_rule_17 = fees_df[fees_df['ID'] == 17]
2224
               if not fee_rule_17.empty:
                   # Use .iloc[0] to get the rule as a Series
fee_rule_17 = fee_rule_17.iloc[0]
2225
2226
                    \sharp Step 4: Filter the transactions that match all conditions of the fee rule.
2227
                    # The manual states that a null or empty list in a rule field means it applies to all values.
                   mask = pd.Series(True, index=full_transaction_data.index)
2228
                    # Match card_scheme (string)
                   if fee_rule_17['card_scheme']:
                        mask &= (full_transaction_data['card_scheme'] == fee rule 17['card scheme'])
2231
                   # Match account_type (list)
                   if fee_rule_17['account_type']: # An empty list evaluates to False, correctly skipping the filter
2232
                        mask &= full_transaction_data['account_type'].isin(fee_rule_17['account_type'])
2233
                    # Match merchant_category_code (list)
if fee_rule_17['merchant_category_code']:
2234
                        mask &=
2235
                        full_transaction_data['merchant_category_code'].isin(fee_rule_17['merchant_category_code'])
                   # Match aci (list)
2237
                   if fee_rule_17['aci']:
                        mask &= full_transaction_data['aci'].isin(fee_rule_17['aci'])
2238
2239
                    # Match is_credit (boolean, but can be null)
                    if fee_rule_17['is_credit'] is not None:
2240
                        mask &= (full_transaction_data['is_credit'] == fee_rule_17['is_credit'])
2241
                   # Apply the combined mask to get the affected transactions
                   affected transactions df = full transaction data[mask]
2242
2243
                   # Step 5: Calculate the original and new fees for the affected transactions.
                   if not affected_transactions_df.empty:
2244
                        total_affected_volume = affected_transactions_df['eur_amount'].sum()
                        num_affected_transactions = len(affected_transactions_df)
2245
                        # Get fee parameters from the rule
                        original_rate = fee_rule_17['rate']
2247
                        fixed_amount = fee_rule_17['fixed_amount']
                        new_rate = 99
2248
                        fee_divisor = 10000
2249
                       \# Calculate total original fee using the formula: fee = fixed_amount + rate * value / 10000 total_fixed_component = fixed_amount * num_affected_transactions
                        total_original_rate_component = (original_rate * total_affected_volume) / fee_divisor
2251
                        total_original_fee = total_fixed_component + total_original_rate_component
2252
                        # Calculate total new fee with the modified rate
                        total_new_rate_component = (new_rate * total_affected_volume) / fee divisor
                        total_new_fee = total_fixed_component + total_new_rate_component
2254
                        # The delta is the difference between the new and original total fees
2255
                        fee_delta = total_new_fee - total_original_fee
2256
                 Final Output
2257
           print(f"{fee_delta:.14f}")
2258
```

Listing 42: **Final solution code.** 

2262 2263

```
f'''2.6772720000000'''
```

Listing 43: Generated answer by DS-STAR. Obtained by executing the above final solution code.

# G PROMPTS

2268

22692270

2271

2272

2273 2274

2275

2276

2277

2278

2279

2280

2282

2283

2284

2285 2286

2287 2288

2289

2290

2291

2292

2293

2294

2295

2296 2297

2298

2299

2300

2301

23022303

23042305

2306

2307

2308

2309

2311

2312 2313

2314

2315

2316

2317

2318

2319

2320 2321

### G.1 ANALYZER AGENT

```
f'''You are an expert data analysist.
Generate a Python code that loads and describes the content of
# Requirement
- The file can both unstructured or structured data.
- If there are too many structured data, print out just few examples.
- Print out essential informations. For example, print out all the

→ column names.

- The Python code should print out the content of {filename}.
- The code should be a single-file Python program that is self-contained
\rightarrow and can be executed as-is.
- Your response should only contain a single code block.
- Important: You should not include dummy contents since we will debug
→ if error occurs.
- Do not use try: and except: to prevent error. I will debug it
→ later.'''
```

Listing 44: Prompt used to generate Python scripts that describe data files.

DS-STAR begins by creating a description of each data file. Specifically, it utilizes an analyzer agent  $\mathcal{A}_{analyzer}$  that takes a file name (e.g., 'payment.csv') as input to generate a Python script  $s_{desc}$ , using the above prompt. The script is then executed to summarize the dataset's core information. An example of this process, including the generated script and its output, is provided in Appendix D.

## G.2 PLANNER AGENT

```
f'''You are an expert data analysist.
In order to answer factoid questions based on the given data, you have

→ to first plan effectively.

# Ouestion
{question}
# Given data: {filenames}
{filenames #1}
{summaries #1}
{filenames #N}
{summaries #N}
# Your task
- Suggest your very first step to answer the question above.
- Your first step does not need to be sufficient to answer the question.
- Just propose a very simple initial step, which can act as a good

→ starting point to answer the question.

- Your response should only contain an initial step.'''
```

Listing 45: Prompt used to generate an initial plan.

DS-STAR begins the solution generation process after generating analytic descriptions of N data files. First, a planner agent  $\mathcal{A}_{\texttt{planner}}$  generates an *initial* high-level executable step  $p_0$  using the above prompt. Here, the query q and obtained data descriptions are used as inputs.

```
f'''You are an expert data analysist.
In order to answer factoid questions based on the given data, you have
   → to first plan effectively.
Your task is to suggest next plan to do to answer the question.
# Question
{question}
```

```
2322
2323
       # Given data: {filenames}
2324
       {filenames #1}
2325
       {summaries #1}
2326
       {filenames #N}
2327
       {summaries #N}
2328
2329
       # Current plans
2330
       1. {Step 1}
2331
       k. {Step k}
2332
2333
       # Obtained results from the current plans:
2334
       {result}
2335
       # Your task
2336
       - Suggest your next step to answer the question above.
2337
       - Your next step does not need to be sufficient to answer the question,
2338
       → but if it requires only final simple last step you may suggest it.
2339
       - Just propose a very simple next step, which can act as a good
2340
        → intermediate point to answer the question.
       - Of course your response can be a plan which could directly answer the
2341
       \rightarrow question.
2342
       - Your response should only contain an next step without any
2343

→ explanation.'''

2344
```

Listing 46: Prompt used to generate a plan after the initialization step.

After the initial plan is obtained, DS-STAR's planner agent generates a subsequent step using the above prompt. Unlike the initialization round, planner agent  $\mathcal{A}_{analyzer}$  additionally uses the intermediate step-by-step plans, and the obtained results after execution of its implementation.

# G.3 CODER AGENT

```
f'''# Given data: {filenames}
{filenames #1}
{summaries #1}
...
{filenames #N}
{summaries #N}

# Plan
{plan}

# Your task
- Implement the plan with the given data.
- Your response should be a single markdown Python code (wrapped in ```).
- There should be no additional headings or text in your response.'''
```

Listing 47: Prompt used to generate a implementation of an initial plan in a Python code.

The initial step for the plan is implemented as a code script  $s_0$  by a coder agent  $\mathcal{A}_{\text{coder}}$ , using the above prompt. Here, analytic descriptions of data files are additionally used to provide some essential information like column names.

```
f'''You are an expert data analysist.
Your task is to implement the next plan with the given data.

# Given data: {filenames}
{filenames #1}
{summaries #1}
...
{filenames #N}
```

```
2376
       {summaries #N}
2377
2378
       # Base code
        ```python
2379
       {base_code}
2380
2381
2382
       # Previous plans
2383
       1. {Step 1}
2384
       k. {Step k}
2385
2386
       # Current plan to implement
2387
       {Step k+1}
2388
       # Your task
2389
       - Implement the current plan with the given data.
2390
         The implementation should be done based on the base code.
2391
       - The base code is an implementation of the previous plans.
       - Your response should be a single markdown Python code (wrapped in

→ ```).
        - There should be no additional headings or text in your response.'''
2394
2395
```

Listing 48: Prompt used to generate a implementation of a plan after the initialization round.

For every round beyond the initial one, DS-STAR's coder agent  $\mathcal{A}_{\text{coder}}$  converts the current plan step  $p_{k+1}$  into a Python script using the above prompt. Unlike the initial round, the agent is provided with the intermediate solution code that implements all preceding step  $p_0, \dots, p_k$ . This allows the agent to incrementally build upon the existing code, simplifying the implementation of each subsequent step.

## G.4 VERIFIER AGENT

23962397

2398

2399

2400

2401

24022403

24262427

2428

2429

```
2404
       f'''You are an expert data analysist.
2405
       Your task is to check whether the current plan and its code
2406
        → implementation is enough to answer the question.
2407
        # Plan
2408
       1. {Step 1}
2409
2410
       k. {Step k}
2411
        # Code
2412
        ```python
2413
        {code}
2414
2415
        # Execution result of code
2416
        {result}
2417
2418
        # Question
2419
       {question}
2420
        # Your task
2421
        - Verify whether the current plan and its code implementation is enough
2422
        \rightarrow to answer the question.
2423
        - Your response should be one of 'Yes' or 'No'.
2424
       - If it is enough to answer the question, please answer 'Yes'.
        - Otherwise, please answer 'No'.'''
2425
```

Listing 49: Prompt used to verify the current plan.

We introduce a verifier agent  $A_{\text{verifier}}$ . At any given round, this agent evaluates the state of the solution, *i.e.*, whether the plan in sufficient or insufficient to solve the problem. The evaluation is

based on the cumulative plan, the user's query, the solution code, which is an implementation of the cumulative plan, and its execution result, using the above prompt.

#### G.5 ROUTER AGENT

2430

2431

2432 2433

2457

2458 2459

2460

2461

2462

```
2434
       f'''You are an expert data analysist.
       Since current plan is insufficient to answer the question, your task is
2435
       \rightarrow to decide how to refine the plan to answer the question.
2436
2437
       # Ouestion
2438
       {question}
2439
       # Given data: {filenames}
2440
       {filenames #1}
2441
       {summaries #1}
2442
2443
       {filenames #N}
2444
       {summaries #N}
2445
       # Current plans
2446
       1. {Step 1}
2447
2448
       k. {Step k}
2449
       # Obtained results from the current plans:
2450
       {result}
2451
2452
       # Your task
2453
       - If you think one of the steps of current plans is wrong, answer among
2454
       → the following options: Step 1, Step 2, ..., Step K.
       - If you think we should perform new NEXT step, answer as 'Add Step'.
2455
       - Your response should only be Step 1 - Step K or Add Step.'''
2456
```

Listing 50: Prompt used for the router agent which determines how to refine the plan.

If the verifier agent  $A_{verifier}$  determines that the current plan is insufficient to solve the user's query, DS-STAR must decide how to proceed. To this end, DS-STAR employs a router agent  $A_{router}$  which decides whether to append a new step or to correct an existing one, which leverages the above prompt.

## G.6 FINALYZER AGENT

```
2463
       f'''You are an expert data analysist.
2464
       You will answer factoid question by loading and referencing the
2465

→ files/documents listed below.

2466
       You also have a reference code.
2467
       Your task is to make solution code to print out the answer of the
2468
        \hookrightarrow question following the given guideline.
2469
        # Given data: {filenames}
2470
       {filenames #1}
2471
       {summaries #1}
2472
2473
        {filenames #N}
        {summaries #N}
2474
2475
        # Reference code
2476
       ```python{code}
2477
2478
        # Execution result of reference code
2479
       {result}
2480
2481
        # Question
2482
       {question}
2483
       # Guidelines
```

```
2484
       {quidelines}
2485
2486
       # Your task
2487
       - Modify the solution code to print out answer to follow the give

→ guidelines.

2488
       - If the answer can be obtained from the execution result of the
2489

ightarrow reference code, just generate a Python code that prints out the
2490

→ desired answer.

2491
       - The code should be a single-file Python program that is self-contained
2492
           and can be executed as-is.
       - Your response should only contain a single code block.
2493
       - Do not include dummy contents since we will debug if error occurs.
2494
       - Do not use try: and except: to prevent error. I will debug it later.
2495
       - All files/documents are in `data/` directory.'''
2496
```

Listing 51: Prompt used to generate a final solution, which outputs the answer with the desired format.

To ensure properly formatted output, DS-STAR employs a finalyzer agent  $A_{\text{finalyzer}}$ , which takes formatting guidelines, if exist, and generates the final solution code.

## G.7 Debugging agent

```
f'''# Error report
{bug}

# Your task
- Remove all unnecessary parts of the above error report.
- We are now running {filename}.py. Do not remove where the error
- occurred.'''
```

Listing 52: Prompt used to summarize the traceback of the error.

When debugging, we first summarize the obtained traceback of error using the above prompt, since the traceback may be too lengthy.

```
f'''# Code with an error:
   ``python
{code}

# Error:
{bug}

# Your task
- Please revise the code to fix the error.
- Provide the improved, self-contained Python script again.
- There should be no additional headings or text in your response.
- Do not include dummy contents since we will debug if error occurs.
- All files/documents are in `data/` directory.'''
```

Listing 53: Prompt used for debugging when analyzing data files.

When generating  $\{d_i\}_{i=1}^N$  using  $s_{\texttt{desc}}$  obtained from  $\mathcal{A}_{\texttt{analyzer}}$ , our debugging agent iteratively update the script using only the summarized traceback, using the above prompt.

```
2538
       {code}
2539
2540
2541
       # Error:
       {bug}
2542
2543
       # Your task
2544
       - Please revise the code to fix the error.
2545
       - Provide the improved, self-contained Python script again.
2546
       - Note that you only have {filenames} available.
       - There should be no additional headings or text in your response.
2547
       - Do not include dummy contents since we will debug if error occurs.
2548
       - All files/documents are in `data/` directory.'''
2549
```

Listing 54: Prompt used for debugging when generating a solution code.

Once DS-STAR obtains  $\{d_i\}_{i=1}^N$ , our debugging agent utilizes such information with the above prompt when generating a solution Python script.

# H LARGE LANGUAGE MODEL USAGE FOR WRITING

In this paper, Gemini was utilized as a general-purpose writing tool. When draft text was provided to Gemini, it improved the writing (*e.g.*, by correcting grammatical errors). The generated text was then manually revised by the authors.