# Causal AI Scientist: Facilitating Causal Data Science with Large Language Models

# **Anonymous Author(s)**

Affiliation Address email

## **Abstract**

Large Language Models (LLMs) have recently been used to automate the causal inference process by overcoming the expertise barrier. However, existing LLMpowered approaches for causal effect estimation often require human users to manually specify variables and methods, and those that do not require manual specification support only a limited set of causal effect measures. To address these limitations, we present Causal AI Scientist (CAIS), an LLM-augmented causal tool with self-correction capabilities. Specifically, given a natural language query and a dataset along with its description, CAIS uses LLMs to understand the user query and dataset, and then selects a method based on a decision tree approach. It then executes the selected method, applies a validation feedback loop for self-correction, and uses the results to answer the input question, enabling fully autonomous causal analysis. Extensive experiments across diverse queries curated from textbooks, synthetic data, and real-world datasets demonstrate CAIS's ability to produce precise causal effect estimates through improved method selection and self-corrections, while reducing runtime errors. We believe CAIS will serve as a strong foundation for enabling fully automated causal inference with LLMs.

#### 1 Introduction

2

3

4

6

8

9

10

11

12

13

14 15

16

Causal inference [Pearl, 2009, Imbens and Rubin, 2015] aims to quantify the effect of a treatment or intervention on an outcome of interest. Understanding cause–effect relationships is central to evidence-based decision-making in fields such as social science [Imbens, 2024], public health [Glass 20 et al., 2013], and biomedicine [Kleinberg and Hripcsak, 2011]. Conducting rigorous causal analysis 21 typically requires methodological expertise, from selecting valid causal effect measures (i.e., esti-22 mands) to choosing appropriate statistical methods, which limits accessibility for non-experts and 23 poses significant challenges to fully automating the causal inference pipeline. For instance, a policy 24 analyst with education and wage data may wish to estimate the effect of a job training program on 25 earnings but, without causal identification knowledge, could reach invalid conclusions. 26

Recently, Large Language Models (LLMs) have emerged as a solution to overcoming the expertise barrier, as they can automate parts of the causal inference process using their extensive knowledge across various domains [Kiciman et al., 2024]. Jiang et al. [2024a] have developed a specialized foundational model, LLM4Causal, for causal inference and causal graph learning. Similarly, more recent approaches have developed causal agents that leverage general-purpose foundational models like GPT [OpenAI et al., 2024] to enable end-to-end performance of causal inference and learning tasks [Wang et al., 2025, Han et al., 2024].

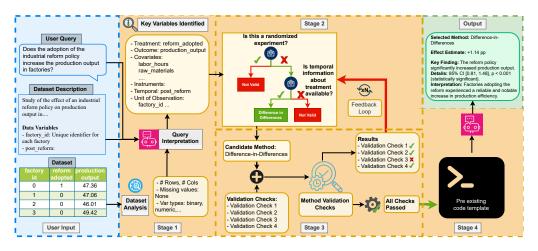


Figure 1: **CAIS workflow.** The user provides an input dataset (CSV file), its description, and an associated causal query. Guided by a decision tree and a backbone LLM, CAIS selects an appropriate estimation method, executes the code, and returns the estimated causal effect along with a natural-language interpretation.

However, existing approaches face three main limitations: (i) fine-tuned models are often narrow in scope, supporting only a limited set of causal effect measures and excluding widely used methods in applied research; (ii) general-purpose tools are primarily evaluated on causal discovery tasks, leaving their causal estimation capabilities untested; and (iii) those that are tested for causal effect estimation are evaluated on examples where users specifically mention the causal effect to be computed, which assumes prior knowledge of causal inference. These limitations are insufficient for enabling the automation of causal inference that is both comprehensive and accessible to others.

To address these limitations, we present Causal AI Scientist (CAIS), an end-to-end LLM-augmented causal tool for generating causality-driven answers to natural language queries. Given a dataset, its description, and a query, CAIS frames the task as a causal inference problem, automatically selects an appropriate method, estimates the causal effect, and interprets the result in context, as outlined in Figure 1. To select the right method, CAIS uses a structured decision tree that breaks down the selection process into focused steps. At each node, it prompts LLMs to evaluate specific features of the dataset or query, such as identifying the treatment, outcome, or instrument. This step-by-step approach reduces errors that commonly arise when relying entirely on LLMs for model selection. Additionally, CAIS performs diagnostic checks and incorporates a feedback loop to self-correct potential errors before producing a final answer.

We evaluate CAIS on CauSciBench [Acharya et al., 2025], a real-world benchmark designed to evaluate the ability of LLMs to perform causal analysis on tabular datasets. CauSciBench consists of causal queries curated from real-world empirical studies, textbook-based examples from QRData (a standard benchmark dataset for causal inference developed by Liu et al. [2024a]), and simulated scenarios. Experiments on its three subsets, real-world studies, textbook examples, and synthetic datasets, show that CAIS outperforms other baselines in selecting the appropriate causal inference method and estimating accurate causal effect values.

58 In summary, our contributions are:

- We propose Causal AI Scientist (CAIS), the first fully autonomous tool designed specifically for causal inference.
- We introduce a structured, rule-based decision tree that guides the method selection process.
  - We perform rigorous evaluations across multiple datasets, and conduct ablation studies to demonstrate the effectiveness of CAIS.

## 4 2 Related Work

LLMs and Causal Inference LLMs have been applied to causal inference with text data [Dhawan 65 et al., 2024, Lin et al., 2023, Imai and Nakamura, 2024, Veljanovski and Wood-Doughty, 2024]. Recent research has also explored their use for causal effect estimation in tabular datasets [Liu et al., 67 2024b, Chen et al., 2025]. However, these approaches often require users to specify the estimation 68 method or variables. Jiang et al. [2024b] introduced a fine-tuned model for causal discovery and effect 69 estimation, but it does not support methods like Instrumental Variables and Difference-in-Differences. 70 Causal-Copilot [Wang et al., 2025] expands the range of methods but was primarily evaluated on 71 causal discovery, not causal inference. Another approach builds causal graphs with LLMs [Kiciman 72 et al., 2024, Han et al., 2024], but this is limited to graph-based estimation techniques. In contrast, 73 our model automates variable and method selection using a decision tree and adds self-correction, 74 creating a fully automatic causal inference pipeline for tabular datasets. 75

**LLM-powered data analysis** Several works have studied the code generation capabilities of 76 77 LLMs for data science tasks, including machine learning, statistical analysis, data manipulation, and visualization [Huang et al., 2022, Lai et al., 2023, Cheng et al., 2023, Nejjar et al., 2024, Jansen 78 et al., 2023]. However, these approaches require users to provide specific instructions. Wu et al. 79 [2024] extend this line of work by enabling LLM-powered tools to perform statistical reasoning and 80 generate solutions to natural language questions. However, these do not involve causal methods. A promising direction for end-to-end analysis is the development of LLM-powered agents. Most 82 of these tools are geared toward machine learning tasks [Zhang et al., 2023, 2024, Huang et al., 83 2024] or data science tasks involving both machine learning and statistical methods [Guo et al., 2024, Hong et al., 2024]. The capabilities of these tools have been enhanced through case-based reasoning [Guo et al., 2024], hierarchical decomposition [Hong et al., 2024], and interactive tools [Wu et al., 86 2023]. However, these agents do not focus on causality-based analysis, which requires different 87 methodological considerations. 88

#### 3 Problem Formulation

92

93

94

90 We are tasked with developing a system that can automatically perform causal effect estimation. The
91 system receives three inputs:

- A tabular dataset,  $\mathcal{D}$ , containing observations for multiple units.
- Metadata describing the variables and the data collection process.
- A natural-language query, q, posing a causal question about the relationships within the data.

The objective is to interpret the query q in the context of the dataset  $\mathcal{D}$  and its metadata to produce an estimate of the causal effect of interest.

#### 97 3.1 Causal Model and Estimand

To formalize this task, we adopt the potential outcomes framework [Rubin, 2005]. The system must first parse the query q and the metadata to identify the key variables: the **treatment** (T), the **outcome** (Y), and a set of covariates (X). Each unit i has two potential outcomes:  $Y_i(1)$ , the outcome if the unit received the treatment, and  $Y_i(0)$ , the outcome if it did not. Similarly,  $T_i$  denotes whether unit i belongs to the treatment  $(T_i = 1)$  or the control  $(T_i = 0)$  group.

One of the causal estimands of interest is the Average Treatment Effect (ATE), defined as the expected difference between these potential outcomes across the population:

$$\tau_{ATE} = \mathbb{E}[Y(1) - Y(0)]$$

However, the fundamental problem of causal inference is that for any given unit, we can only observe one of its potential outcomes [Holland, 1986]. The observed outcome is  $Y_i = Y_i(T_i)$ . This means the ATE cannot be calculated directly.

Example: Job Training Program To make this concrete, consider a dataset from a job-training program (e.g., the Lalonde dataset [LaLonde, 1986]) and the query, "Does the training program boost earnings?". Here, the treatment T is program participation, the outcome Y is earnings, and

covariates *X* could include education, age, and prior income. The ATE would represent the average boost in earnings if everyone in the population participated in the program versus if no one did.

#### 3.2 Identification: From Randomized to Observational Data

111

137

144

147

The process of connecting the unobservable ATE to our observed data is called **identification**. This requires making assumptions.

The simplest case is a Randomized Controlled Trial (RCT), where individuals are randomly assigned to treatment (T=1) or control (T=0). This randomization ensures that, on average, the two groups are identical before treatment. This satisfies the **ignorability** assumption  $((Y(1), Y(0)) \perp \!\!\! \perp T)$ , meaning the treatment assignment is independent of the potential outcomes. In an RCT, the ATE can be simply estimated by the difference in the average outcomes of the two groups:

$$\hat{\tau}_{ATE} = \frac{1}{N_1} \sum_{i:T_i=1} Y_i - \frac{1}{N_0} \sum_{i:T_i=0} Y_i.$$
 (1)

However, most data is **observational**, not from an RCT. In observational data, ignorability is often violated. For instance, more motivated individuals (who might have higher earnings potential anyway) may be more likely to sign up for a job training program. This motivation is a **confounder**, a variable that affects both treatment assignment and the outcome.

To handle confounders, we rely on a stronger, untestable assumption: **conditional ignorability**. This states that the treatment assignment is independent of the potential outcomes when conditioned on all the common causes (X) of T and Y. Mathematically, it can be formulated as follows:

$$(Y(1), Y(0)) \perp \!\!\!\perp T \mid X \tag{2}$$

Along with the **Stable Unit Treatment Value Assumption** (**SUTVA**) (no interference between units and consistency of treatment), conditional ignorability allows us to identify the ATE by adjusting for the covariates X:

$$\tau_{ATE} = \mathbb{E}_X[\mathbb{E}[Y \mid T = 1, X] - \mathbb{E}[Y \mid T = 0, X]]$$

Estimation via Statistical Models To compute the quantity above, we use structural causal models [Pearl, 2009] to describe the relationships between Y, T, and X. A common approach is a linear structural model, which can be expressed as below:

$$Y = \alpha + X^T \beta + \tau T + \epsilon, \tag{3}$$

where  $\alpha$  is the intercept,  $\beta$  is a vector of coefficients for the covariates  $X, \tau$  is the treatment effect parameter, and  $\epsilon$  is the unobserved error term. In this model, the coefficient  $\tau$  directly corresponds to the ATE, as it represents the change in Y for a one-unit change in T after adjusting for X. We can then use an estimation method, such as linear regression, to find a sample-based estimate,  $\hat{\tau}$ , from the data.

#### 3.3 The Role of the LLMs: Supplying Domain Expertise

The entire causal inference process hinges on the validity of the conditional ignorability assumption, which cannot be verified from data alone. Justifying this assumption requires **domain knowledge** to argue that the set of measured covariates X is sufficient to account for all major confounders. This is precisely where the LLMs come in. The LLMs are designed to act as a proxy for a human domain expert. By leveraging their vast knowledge base, the LLMs can analyze the dataset's metadata and the user's query to:

- Propose a plausible causal relationship between the variables.
- Identify the most likely confounders that must be included in the set X.
- Justify the conditional ignorability assumption, thereby enabling a principled estimation of the causal effect.

# 8 4 Methodology: CAIS

- CAIS consists of four successive methodological stages, each consisting of one or more micro-tools.
  Every stage combines logic from established causal inference principles with LLM reasoning, which
  is selectively applied to sub-tasks that require human-like judgment. The overall framework is
- depicted in Figure 1. Detailed descriptions of the specific micro-tools used in each stage are provided in Appendix D.
- Outline Specifically, our method consists of the following stages:
- Stage 1 (Data Preprocessing & Query Decomposition) The process begins by analyzing the dataset to identify key components, such as control and target variables (Section 4.1).
- Stage 2 (Method Selection) A rule-based decision tree is used to select a valid method for causal effect estimation based on the identified components (Section 4.2).
- **Stage 3 (Validation)** The standard assumptions for the selected method are then validated. If any assumption check fails, the system backtracks to Stage 2 to find an alternative method, creating a validation loop (Section 4.3).
- Stage 4 (Method Execution & Interpretation) Finally, once all checks pass, the chosen method is executed using predefined templates, and the final result is returned with an interpretation (Section 4.4).

#### 4.1 Stage 1: Dataset Preprocessing & Query Decomposition

In this stage, the CAIS uses an LLM to analyze the dataset and formalize the user's causal question. 166 Guided by a structured prompt (see Appendix E.1 for details), in this stage we perform two sequential 167 tasks. First, we profile the data by generating a summary of column types, missing values, and 168 statistical distributions. Using this summary as context, we use an LLM to decompose the user's 169 query. It interprets the natural language request to identify and categorize the essential columns for 170 the analysis. This involves designating the treatment, outcome, and confounder variables, as well 171 as scanning for specialized variables—such as instrumental, running, or time-series variables—that 172 enable specific causal methods. The final output is a structured definition of the causal problem, 173 which serves as the input for Stage 2. 174

# 175 4.2 Stage 2: Method Selection

165

189

In this stage, CAIS selects a suitable causal inference method using the structured specification 176 produced in Stage 1. The selection is performed through a rule-based decision tree that encodes 177 standard design logic from the causal inference literature. At each branch, the system evaluates a 178 specific property of the problem, such as whether treatment is randomly assigned, whether time 179 and unit identifiers are present, or whether an instrumental or running variable is available. This 180 information is already extracted in Stage 1 with deterministic rules and LLM assistance wherever 181 necessary, ensuring transparency and reproducibility. The tree then routes the problem to one of a 182 small number of valid methods, such as difference-in-means, ordinary least squares, difference-in-183 differences, etc. 184

By breaking the selection process into explicit and verifiable steps, the decision tree ensures both accuracy and interpretability, avoiding the opacity of direct method selection by an LLM. The resulting method choice, along with its assumptions, is then passed to Stage 3 for validation. We provide a detailed explanation of the decision tree for method selection in Appendix C.

# 4.3 Stage 3: Validation

This stage serves as a crucial validation and feedback mechanism, acting as a safeguard against errors from the initial analysis and method selection. CAIS performs the standard **statistical assumption**checks required for the selected method, such as the parallel trends test for Difference-in-Differences or computing the F-statistic for an Instrumental Variable (IV) analysis. If failure occurs for any of these assumptions, the system initiates a feedback loop back to Stage 2. Information from the failed validation attempt is incorporated into the context, allowing us to skip past the previously selected

Collection	# Queries	# CSV	Median Obs.	Median Cols.
QRData	39	35	1209	19.0
RealPapers	29	14	1720	17.5
Synthetic	45	45	428	7.0

Table 1: Statistics of the CauSciBench dataset

node from the decision tree and move on to the next plausible candidate. We provide a detailed qualitative example of the validation loop in Appendix F.

#### 4.4 Stage 4: Method Execution & Interpretation

Once a method successfully passes all validation checks in Stage 3, CAIS proceeds to **execution**. This stage utilizes predefined code templates with placeholders for the key variables (e.g., treatment, outcome) identified in Stage 1. This template-based strategy was chosen over LLM-powered code generation to maximize reliability and efficiency. While generating code from scratch can be flexible, it increases the risk of implementation errors and can be slow and costly due to the need for iterative debugging [Chen et al., 2025]. Our approach minimizes these risks, as its core logic is pre-verified. After the method is executed, the final step is **interpretation**. The LLM is prompted to synthesize the numerical output—such as the causal estimate and its statistical significance—into a natural language explanation that directly addresses the user's original query. Crucially, this interpretation is presented alongside important caveats, including the results of the validation checks from Stage 3 and a clear statement of the assumptions and limitations of the chosen method. This ensures the user understands the full context and reliability of the final estimate.

# 5 Experimental Setup

Dataset To evaluate CAIS, we use CauSciBench [Acharya et al., 2025], a comprehensive benchmark for assessing LLMs on causal-estimation tasks. CauSciBench consists of 113 causal queries drawn from three sources: (1) QRData [Liu et al., 2024b], a benchmark dataset primarily based on causal inference textbooks; (2) published papers across multiple disciplines; and (3) synthetic datasets with known ground-truth causal effects. This collection covers a wide range of evaluation scenarios and methods used in practice. Summary statistics of CauSciBench are presented in Table 1.

**Baselines** Given the lack of LLM-based tools for fully automated causal inference, we compare our method against three strong prompting strategies that represent the state of the art in LLM-assisted data analysis. ReAct prompting [Yao et al., 2023] guides the model through iterative thought–action–observation cycles, which prior work shows to be highly effective for causal inference [Liu et al., 2024b]. Program of Thoughts (PoT) prompting [Chen et al., 2022] instead asks the model to produce a single, complete program that handles the entire analysis from data loading to reporting results. Finally, a Veridical Data Science–inspired prompt [Yu, 2020] emphasizes stability by requiring the model to reflect on and critique its own methodological choices before finalizing outputs. Examples of these prompts are provided in Appendix G.

**Implementation Details** For estimating causal effects, we use the DoWhy [Sharma and Kiciman, 2020, Blöbaum et al., 2024] and statsmodels [Seabold and Perktold, 2010] Python libraries. Our experiments use several LLMs as the reasoning backbone, including GPT-40, Llama-3.3-70B-Instruct, and Gemini 2.5 Pro. All models were accessed via their respective APIs. To ensure reproducibility for all experiments, we used greedy decoding by setting the temperature parameter to 0.

**Evaluation Metrics** We evaluate our pipeline using the following metrics.

• Method Selection Accuracy (MSA): Percentage of queries where the selected method  $\hat{m}_i$  matches the reference method  $(m_i)$ .

$$MSA = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}[\hat{m}_i = m_i] \times 100.$$
 (4)

• Mean Relative Error (MRE): Average relative error between predicted causal effects  $(\hat{\tau}_i)$  and reference values  $(\tau_i)$ .

$$MRE = \frac{1}{N} \sum_{i=1}^{N} \min\left(\frac{|\hat{\tau}_i - \tau_i|}{|\tau_i|}, 1\right) \times 100\%$$
 (5)

To reduce sensitivity to outliers, relative error is capped at 100% per query.

In the above formulas, N denotes the total number of queries in the evaluation set.

#### 6 Results

We evaluate the performance of CAIS against the baseline prompting strategies (Section 5). We then analyze the effectiveness of its key components, the decision tree and the validation feedback loop, through ablation studies (Section 6.1). We also provide a detailed comparative performance analysis in Appendix A and detailed error analysis B. Table 2 presents the main results, comparing CAIS with the baselines on three subsets of CauSciBench.

Method Selection Accuracy (↑)			Mean Relative Error (↓)							
Method	GPT-40	GPT-40-mini	o3-mini	Gemini 2.5 Pro	Llama 3.3 70B	GPT-40	GPT-40-mini	o3-mini	Gemini 2.5 Pro	Llama 3.3 70B
					Textbook De	ata				
ReAct	55.0	55.2	21.8	62.2	34.4	43.2	33.9	44.7	43.2	43.9
PoT	41.0	54.3	30.7	50.0	53.8	32.6	33.6	30.7	35.8	31.5
Veridical	60.5	41.0	61.5	59.0	46.1	40.7	42.2	27.6	37.8	55.4
CAIS	74.4	74.3	94.1	81.2	81.8	31.6	55.9	43.1	41.2	54.2
					Synthetic De	ata				
ReAct	51.2	41.9	46.7	48.2	55.8	27.9	21.2	21.0	20.2	21.3
PoT	53.3	37.7	42.2	53.2	47.6	19.9	37.7	42.2	24.0	21.1
Veridical	79.0	43.4	66.6	58.5	50.0	27.7	25.7	20.2	26.5	33.3
CAIS	76.9	75.9	73.3	75.6	79.5	17.4	16.2	20.0	18.5	50.0
					Real Data	ı				
ReAct	69.5	51.8	57.1	55.0	44.4	43.1	52.3	43.2	38.1	52.6
PoT	57.7	54.6	33.3	42.2	53.8	54.7	55.6	46.3	42.0	53.8
Veridical	48.0	28.0	59.2	53.2	24.0	53.6	54.4	41.2	39.0	52.8
CAIS	69.2	65.2	76.9	78.3	73.0	47.5	54.6	39.7	32.0	37.4

Table 2: Performance of CAIS and baseline prompting strategies across all datasets and LLMs. Results are reported for both **Method Selection Accuracy (MSA, higher is better)** and **Mean Relative Error (MRE, lower is better)**. Dataset blocks correspond to Textbook, Synthetic, and Real-world settings. Bold entries indicate the best value in each row. CAIS consistently outperforms baselines on MSA while maintaining competitive MRE.

CAIS shows superior method selection capabilities. Across all three datasets and models, we observe that CAIS outperforms the baselines in MSA, with significant margins in nearly all cases. For example, on the textbook dataset, o3-mini achieves an MSA of 94.1%, which is 32.6 percentage points higher than the second-best baseline. On average, CAIS improves MSA over the best-performing baseline per LLM by +22.18 points on the textbook data, +15.58 points on the synthetic data, and +14.10 points on the real data. These results demonstrate the effectiveness of our decision-tree-based method selection and validation feedback loop compared to relying solely on LLM reasoning.

**CAIS** achieves competitive causal estimation accuracy. Performance trends in MRE are more nuanced. On the synthetic and real datasets, our model performs competitively, achieving either the lowest or second-lowest MRE for most LLMs. However, a notable exception is the Textbook dataset, where our model consistently underperforms.

We attribute this discrepancy to a key design choice intended to prevent incorrect implementation. Our model uses pre-verified code templates and does not retry execution if an error occurs, which can result in a substantial MRE penalty on failure. In contrast, the baseline methods use iterative retries to ensure they always return a numerical estimate, even if the chosen causal method is fundamentally inappropriate for the data.

Study	Method	GPT-40	GPT-40-mini	o3-mini
QR	LLM	48.4	50.6	50.0
	Tree	<b>74.4</b>	<b>74.3</b>	<b>94.1</b>
Real	LLM	48.0	60.8	45.5
	Tree	<b>69.2</b>	<b>65.2</b>	<b>76.9</b>
Synth	LLM	57.5	<b>79.4</b>	57.1
	Tree	<b>76.9</b>	78.9	<b>73.3</b>

Study	Loop	GPT-40	GPT-40-mini	o3-mini
QR	No Yes	<b>80.0</b> 74.4	41.2 <b>74.3</b>	97.1 <b>94.1</b>
Real	No	56.0	33.3	75.0
	Yes	<b>69.2</b>	<b>65.2</b>	<b>76.9</b>
Synth	No	71.4	60.0	77.3
	Yes	<b>76.9</b>	<b>75.9</b>	<b>73.3</b>

<sup>(</sup>a) Decision tree vs. LLM-based method selection.

Table 3: Ablation results: (a) decision tree vs. LLM-based method selection, and (b) validation feedback loop.

#### 6.1 Ablation Studies

We conduct two ablation studies to assess the role of the core components in our pipeline: the *decision tree* and the *validation loop*. For the decision tree ablation, we compare CAIS with a variant that removes the decision tree and prompts an LLM to directly select the method. For the validation loop ablation, we compare CAIS with a variant without the validation loop, which simply uses the initially selected method without further checks.

Decision tree significantly improves method-selection accuracy. The results in Table 3a clearly show that incorporating a structured decision tree for method selection leads to substantially higher MSA compared to relying solely on the LLM's direct judgment. The decision tree decomposes the selection process into a sequence of targeted diagnostic questions, each focused on a specific dataset property (e.g., treatment timing, presence of instruments, covariate balance). This step-by-step approach constrains the LLM's reasoning to smaller, well-defined decisions, reducing the likelihood of overgeneralization or bias toward familiar methods. For example, on QRData, GPT-4o improves from 48.4% to 74.4% and o3-mini from 50% to 94.1%. Overall, the decision-tree-based approach is superior in performance compared to solely relying on LLMs.

The validation loop helps reduce the performance gap between weaker LLMs and stronger LLMs. As shown in Table 3b, the validation feedback loop is a critical component, especially for correcting errors from less capable models. This effect is most evident with GPT-4o-mini, whose accuracy improves dramatically from 41.2% to 74.3% (+33.1 points) on QRData and from 33.3% to 65.2% (+31.9 points) on the Real dataset. The loop provides this significant boost by allowing weaker models to recover from incorrect initial method selections, to which they are more prone.

Conversely, the benefits are less pronounced for stronger models like GPT-40. Since these models are more likely to select the correct method on the first attempt, the validation loop offers only marginal gains and can occasionally lead to a slight decrease in performance.

## 7 Conclusion

In this work, we introduce Causal AI Scientist (CAIS), an end-to-end tool that maps natural language queries and datasets to formal causal inference tasks by automatically selecting appropriate methods and interpreting results. When evaluated across diverse causal inference tasks using three datasets, CAIS consistently outperforms baseline prompting strategies in method selection and achieves competitive performance in causal effect estimation, particularly on structured datasets such as QRData and synthetic examples. These results highlight the value of CAIS's decision-tree-based approach, which decomposes complex reasoning into interpretable steps. This not only improves estimation accuracy but also enhances robustness and transparency—qualities critical for researchers and practitioners in social science, healthcare, and related fields. Moreover, CAIS's strong performance on well-structured datasets suggests that real-world outcomes can be further improved with better data preprocessing, reinforcing its utility as a trustworthy tool for non-experts seeking accessible and interpretable causal analysis.

<sup>(</sup>b) Impact of the validation feedback loop.

#### 8 References

- Sawal Acharya, Terry Jingchen Zhang, Andrew Kim, Anahita Haghighat, Xianlin Sun, Maximilian Mordig, Rahul Babu Shrestha, Furkan Danisman, Clijo Jose, Yahang Qi, Pepijn Cobben,
  Bernhard Schölkopf, Mrinmaya Sachan, and Zhijing Jin. Causcibench: Assessing Ilm causal
  reasoning for scientific research. 2025. URL https://zhijing-jin.com/files/papers/2025\_
  CauSciBench.pdf.
- Patrick Blöbaum, Peter Götz, Kailash Budhathoki, Atalanti A. Mastakouri, and Dominik Janzing.
   Dowhy-gcm: An extension of dowhy for causal inference in graphical causal models. *Journal of Machine Learning Research*, 25(147):1–7, 2024. URL http://jmlr.org/papers/v25/22-1258.
   html.
- David Card. Using geographic variation in college proximity to estimate the return to schooling.
  Working Paper 4483, National Bureau of Economic Research, October 1993. URL http://www.nber.org/papers/w4483.
- David Card and Alan B. Krueger. Minimum wages and employment: A case study of the fast-food industry in new jersey and pennsylvania. *The American Economic Review*, 84(4):772–793, 1994. ISSN 00028282. URL http://www.jstor.org/stable/2118030.
- Christopher Carpenter and Carlos Dobkin. The effect of alcohol consumption on mortality: Regression discontinuity evidence from the minimum drinking age. *Am Econ J Appl Econ*, 1(1):164–182, Jan 2009. ISSN 1945-7782 (Print); 1945-7790 (Electronic); 1945-7790 (Linking). doi: 10.1257/app.1. 1.164.
- Qiang Chen, Tianyang Han, Jin Li, Ye Luo, Yuxiao Wu, Xiaowei Zhang, and Tuo Zhou. Can ai master econometrics? evidence from econometrics ai agent on expert-level tasks, 2025. URL https://arxiv.org/abs/2506.00856.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint* arXiv:2211.12588, 2022.
- Liying Cheng, Xingxuan Li, and Lidong Bing. Is GPT-4 a good data analyst? In *The 2023 Conference*on Empirical Methods in Natural Language Processing, 2023. URL https://openreview.net/
  forum?id=PxEhoPiBB0.
- Nikita Dhawan, Leonardo Cotta, Karen Ullrich, Rahul Krishnan, and Chris J. Maddison. End-to-end causal effect estimation from unstructured natural language data. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=gzQARCgIsI.
- Thomas A. Glass, Steven N. Goodman, Miguel A. Hernán, and Jonathan M. Samet. Causal inference in public health. *Annual review of public health*, 34:61–75, March 2013. ISSN 0163-7525. doi: 10.1146/annurev-publhealth-031811-124606.
- Noah Greifer. Assessing Balance, May 2025. URL https://cran.r-project.org/web/packages/ MatchIt/vignettes/assessing-balance.html. R package MatchIt vignette.
- 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, 2024. URL https://arxiv.org/abs/2402.17453.
- Kairong Han, Kun Kuang, Ziyu Zhao, Junjian Ye, and Fei Wu. Causal agent based on large language model, 2024. URL https://arxiv.org/abs/2408.06849.
- Keisuke Hirano and Guido W. Imbens. The propensity score with continuous treatments. *Applied Bayesian Modeling and Causal Inference from Incomplete-Data Perspectives*, pages 73–84, 2004.
- Paul W. Holland. Statistics and causal inference. *Journal of the American Statistical Association*, 81 (396):945–960, 1986. doi: 10.1080/01621459.1986.10478354. URL https://www.tandfonline.com/doi/abs/10.1080/01621459.1986.10478354.

- Paul W. Holland. Causal inference, path analysis, and recursive structural equations models. *Sociological Methodology*, 18:449–484, 1988. ISSN 00811750, 14679531. URL http: //www.jstor.org/stable/271055.
- Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Chenxing Wei,
  Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Li Zhang, Lingyao Zhang, Min Yang, Mingchen
  Zhuge, Taicheng Guo, Tuo Zhou, Wei Tao, Xiangru Tang, Xiangtao Lu, Xiawu Zheng, Xinbing
  Liang, Yaying Fei, Yuheng Cheng, Zhibin Gou, Zongze Xu, and Chenglin Wu. Data interpreter:
  An Ilm agent for data science, 2024. URL https://arxiv.org/abs/2402.18679.
- Junjie Huang, Chenglong Wang, Jipeng Zhang, Cong Yan, Haotian Cui, Jeevana Priya Inala, Colin
  Clement, and Nan Duan. Execution-based evaluation for data science code generation models.
  In Eduard Dragut, Yunyao Li, Lucian Popa, Slobodan Vucetic, and Shashank Srivastava, editors,

  Proceedings of the Fourth Workshop on Data Science with Human-in-the-Loop (Language Advances), pages 28–36, Abu Dhabi, United Arab Emirates (Hybrid), December 2022. Association
  for Computational Linguistics. URL https://aclanthology.org/2022.dash-1.5/.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents on machine learning experimentation, 2024. URL https://arxiv.org/abs/2310.03302.
- Kosuke Imai and Kentaro Nakamura. Causal representation learning with generative artificial intelligence: Application to texts as treatments, 2024. URL https://arxiv.org/abs/2410. 00903.
- Guido W. Imbens. Instrumental variables: An econometrician's perspective. *Statistical Science*, 29(3): 323–358, 2014. ISSN 08834237, 21688745. URL http://www.jstor.org/stable/43288511.
- Guido W. Imbens. Causal inference in the social sciences. *Annual Review of Statistics and Its Application*, qq:1123–152, 2024.
- Guido W. Imbens and Donald B. Rubin. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction.* Cambridge University Press, 2015.
- Jacqueline A Jansen, Artür Manukyan, Nour Al Khoury, and Altuna Akalin. Leveraging large language models for data analysis automation. *bioRxiv*, 2023. doi: 10.1101/2023.12.11.571140. URL https://www.biorxiv.org/content/early/2023/12/21/2023.12.11.571140.
- Haitao Jiang, Lin Ge, Yuhe Gao, Jianian Wang, and Rui Song. Llm4causal: Democratized causal tools for everyone via large language model, 2024a. URL https://arxiv.org/abs/2312.17122.
- Haitao Jiang, Lin Ge, Yuhe Gao, Jianian Wang, and Rui Song. LLM4causal: Democratized causal tools for everyone via large language model. In *First Conference on Language Modeling*, 2024b. URL https://openreview.net/forum?id=H1Edd5d2JP.
- Emre Kiciman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models: Opening a new frontier for causality. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=mqoxLkX210. Featured Certification.
- Samantha Kleinberg and George Hripcsak. Methodological review: A review of causal inference for biomedical informatics. *J. of Biomedical Informatics*, 44(6):1102–1112, December 2011. ISSN 1532-0464. doi: 10.1016/j.jbi.2011.07.001. URL https://doi.org/10.1016/j.jbi.2011.07.001. 001.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-Tau Yih,
  Daniel Fried, Sida Wang, and Tao Yu. DS-1000: A natural and reliable benchmark for data science
  code generation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,
  Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 18319–18345.
  PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/lai23b.html.
- Robert J LaLonde. Evaluating the Econometric Evaluations of Training Programs with Experimental Data. *American Economic Review*, 76(4):604–620, September 1986. URL https://ideas.repec.org/a/aea/aecrev/v76y1986i4p604-20.html.

- Victoria Lin, Louis-Philippe Morency, and Eli Ben-Michael. Text-transport: Toward learning causal effects of natural language, 2023. URL https://arxiv.org/abs/2310.20697.
- Xiao Liu, Zirui Wu, Xueqing Wu, Pan Lu, Kai-Wei Chang, and Yansong Feng. Are Ilms capable of
   data-based statistical and causal reasoning? benchmarking advanced quantitative reasoning with
   data, 2024a. URL https://arxiv.org/abs/2402.17644.
- Xiao Liu, Zirui Wu, Xueqing Wu, Pan Lu, Kai-Wei Chang, and Yansong Feng. Are LLMs capable of data-based statistical and causal reasoning? benchmarking advanced quantitative reasoning with data. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9215–9235, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.548. URL https://aclanthology.org/2024.findings-acl.548.
- Edward Miguel and Michael Kremer. Worms: Identifying impacts on education and health in the presence of treatment externalities. *Econometrica*, 72(1):159–217, 2004. ISSN 00129682, 14680262. URL http://www.jstor.org/stable/3598853.
- Mohamed Nejjar, Luca Zacharias, Fabian Stiehle, and Ingo Weber. Llms for science: Usage for code
   generation and data analysis. *J. Softw. Evol. Process*, 37(1), September 2024. ISSN 2047-7473.
   doi: 10.1002/smr.2723. URL https://doi.org/10.1002/smr.2723.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni 413 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor 414 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, 415 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny 416 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, 417 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea 418 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 419 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, 420 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, 423 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel 424 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua 425 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike 426 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon 427 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne 428 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo 429 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, 430 431 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, 432 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy 433 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie 434 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, 435 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David 437 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie 438 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, 439 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo 440 Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, 441 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, 442 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, 443 Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, 444 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis 445 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted 446 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel 447 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranay Shyam, Szymon 448 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 449 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie 450

Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,

- 452 Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun
- Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang,
- Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian
- Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren
- Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
- Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
- Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL
- 459 https://arxiv.org/abs/2303.08774.
- Judea Pearl. Causality: Models, Reasoning and Inference. Cambridge University Press, New York,
   2000.
- Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*, 3(none):96 146, 2009.
   doi: 10.1214/09-SS057. URL https://doi.org/10.1214/09-SS057.
- Paul R. Rosenbaum and Donald B. Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983. ISSN 00063444, 14643510. URL http://www.jstor.org/stable/2335942.
- Donald B Rubin. Causal inference using potential outcomes. *Journal of the American Statistical Association*, 100(469):322–331, 2005. doi: 10.1198/016214504000001880. URL https://doi.org/10.1198/016214504000001880.
- Skipper Seabold and Josef Perktold. statsmodels: Econometric and statistical modeling with python.
   In 9th Python in Science Conference, 2010.
- Amit Sharma and Emre Kiciman. Dowhy: An end-to-end library for causal inference. *arXiv preprint arXiv:2011.04216*, 2020.
- Elizabeth A. Stuart. Matching Methods for Causal Inference: A Review and a Look Forward.

  Statistical Science, 25(1):1 21, 2010. doi: 10.1214/09-STS313. URL https://doi.org/10.1214/09-STS313.
- Marko Veljanovski and Zach Wood-Doughty. DoubleLingo: Causal estimation with large language
   models. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 799–807, Mexico City, Mexico,
   June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-short.71.
   URL https://aclanthology.org/2024.naacl-short.71/.
- Xinyue Wang, Kun Zhou, Wenyi Wu, Har Simrat Singh, Fang Nan, Songyao Jin, Aryan Philip, Saloni
   Patnaik, Hou Zhu, Shivam Singh, Parjanya Prashant, Qian Shen, and Biwei Huang. Causal-copilot:
   An autonomous causal analysis agent, 2025. URL https://arxiv.org/abs/2504.13263.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun
   Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and
   Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation, 2023. URL
   https://arxiv.org/abs/2308.08155.
- Xueqing Wu, Rui Zheng, Jingzhen Sha, Te-Lin Wu, Hanyu Zhou, Tang Mohan, Kai-Wei Chang,
   Nanyun Peng, and Haoran Huang. DACO: Towards application-driven and comprehensive data
   analysis via code generation. In *The Thirty-eight Conference on Neural Information Processing* Systems Datasets and Benchmarks Track, 2024. URL https://openreview.net/forum?id=
   NrCPBJSOOC.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
  React: Synergizing reasoning and acting in language models, 2023. URL https://arxiv.org/abs/2210.03629.
- Bin Yu. Veridical data science. In *Proceedings of the 13th International Conference on Web*Search and Data Mining, WSDM '20, page 4–5, New York, NY, USA, 2020. Association for
  Computing Machinery. ISBN 9781450368223. doi: 10.1145/3336191.3372191. URL https:
  //doi.org/10.1145/3336191.3372191.

- Lei Zhang, Yuge Zhang, Kan Ren, Dongsheng Li, and Yuqing Yang. Mlcopilot: Unleashing the power of large language models in solving machine learning tasks, 2024. URL https: //arxiv.org/abs/2304.14979.
- Shujian Zhang, Chengyue Gong, Lemeng Wu, Xingchao Liu, and Mingyuan Zhou. Automl-gpt: Automatic machine learning with gpt, 2023. URL https://arxiv.org/abs/2305.02499.

Metric	Baseline	CAIS	Change (%)
General Statistics			
Total Queries	1551	585	_
Successful Queries	1476	512	_
Total Retries	930	159	_
Retries per Query (%)	59.96	27.18	↓ 54.69
Method Match Rate (%)	52.08	76.20	↑ 46.3
Mean Error (%)	35.38	37.66	↑ 6.4
Error Breakdown (%)			
Execution & Runtime Error	34.39	22.91	↓ 33.4
Method Mismatch	29.77	21.20	↓ 28.8
Data Loading Failure	3.10	0.00	↓ 100.0
Missing Result	0.76	6.84	↑800.0

Table 4: Comparison of performance and error types between baseline and CAIS. Arrow indicates direction of change from baseline to CAIS.

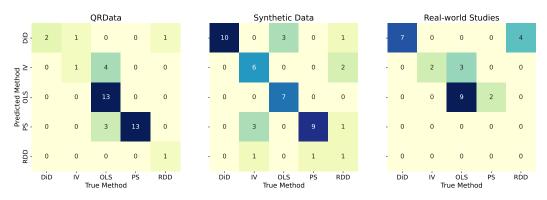


Figure 2: Confusion matrix showing method selection performance of CAIS (GPT-40).

# 507 A Comparative Performance Analysis

Here, we compare the overall performance of CAIS, which uses a structured approach, with baselines that rely on prompting-based LLM code generation, focusing on execution stability, error types, and the frequency of retries.

- **Higher Method Selection Accuracy:** CAIS achieves a 46.3% higher method match rate than the baseline (76.2% vs. 52.08%), indicating more accurate identification of appropriate causal methods.
- Substantial Reduction in Retries: CAIS reduces total retries by 54.6% per query (In CAIS, a retry refers to feedback via validation loop), suggesting more robust and executable outputs due to structured prompt generation and template-based code execution.
- **Improved Execution Stability:** Execution and runtime errors are reduced by 33.4%, and method mismatches decrease by 28.8%, reflecting enhanced reliability in model reasoning and implementation.
- No Data Loading Failures: CAIS handles datasets more reliably with 0% data loading failures compared to 3.1% in the baseline.
- **Trade-offs in Estimation Quality:** While CAIS increases mean error slightly (from 35.38% to 37.66%), this may stem from using more advanced methods rather than defaulting to simple linear regression.

# 525 B Error Analysis

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

548

549 550

551

552

553

554 555

556

557

562

- This section provides a qualitative breakdown of the model's primary failure modes. We analyze common error patterns to identify their root causes and determine which stages of the pipeline are most vulnerable.
  - Incorrect Variable Selection: LLMs frequently misinterpret temporal covariates, such as birth year or quarter indicators, as observation time points. This misinterpretation can erroneously lead to the selection of Difference-in-Differences as the causal inference method. Additionally, LLMs often misidentify treatment and outcome variables, particularly when column names lack clear descriptive labels or contain ambiguous terminology.
  - Wrong Method Selection: As demonstrated in Figure 2, LLMs misclassify Randomized Controlled Trials as Encouragement Designs, leading to the selection of Instrumental Variables instead of linear regression. Similarly, for synthetic datasets, the model fails to identify Instrumental Variables as the optimal method in three instances. This pattern underscores the inherent challenge of selecting valid instruments based solely on data descriptions.
  - Incorrect Data Formats: Implementation errors also stem from inconsistent data formatting.
     Specifically, certain variables are encoded as strings when causal inference packages like DoWhy require numerical inputs, creating compatibility issues that compromise execution.

# C Explanation of decision tree

#### 4 C.1 Notation

- 545 We first define key notation used throughout this guide:
- Y: Outcome (the variable we want to understand or predict)
  - T: Treatment assignment (whether a unit is assigned to treatment; T=0 indicates the control group, T=1 indicates the treatment group)
  - **D**: Treatment uptake (actual receipt of treatment, used in encouragement designs where assignment is not the same as uptake)
  - **Z**: Instrumental variable (a variable used to identify causal effects in the presence of unobserved confounding)
  - M: Mediator (a variable that lies on the causal pathway between treatment and outcome)
  - U: Unobserved confounder (a variable that affects both treatment and outcome but is not measured)
  - X: Covariates (observed variables that may influence treatment or outcome)
    - i: Individual units of analysis (for example people, organizations, or countries)

#### 558 C.2 Randomized Controlled Trials (RCTs)

We begin by determining if the data comes from a Randomized Controlled Trial (RCT). RCTs are the gold standard for causal inference because random assignment eliminates confounding by making treatment and control groups comparable on average.

#### C.2.1 Encouragement Designs

- Within RCTs, we first check if the data comes from an encouragement design. An encouragement design refers to scenarios where treatment assignment is random, but not all assignees accept or take up their assigned treatment, i.e.,  $T_i \neq D_i$  [Holland, 1988].
- A classic example is the deworming experiment by Miguel and Kremer [Miguel and Kremer, 2004].

  Students were randomly assigned to receive deworming drugs, but not all students who were assigned
- actually took the drugs. In such cases, the corresponding relation would be:

T (randomly assigned to get dewormed)  $\rightarrow D$  (actually took the drug)  $\rightarrow Y$  (school enrollment).

To estimate the effect of actual treatment uptake (rather than just assignment), one uses instrumental variable analysis, with the random assignment T serving as an instrument for actual uptake D.

#### 571 C.2.2 Standard RCT Analysis

- 572 If the data is not from an encouragement design, meaning everyone accepts their assigned treatment,
- this is the classic RCT analysis. One can simply compute the difference in means between the
- treatment and control groups.
- 575 In many cases, the dataset might have pretreatment variables (baseline characteristics measured
- before treatment). While one can still compute a simple difference in means, the trend in the literature
- is to include these pretreatment variables as controls in an Ordinary Least Squares regression model.
- This is primarily done to improve the precision of the estimates (reduce the standard errors) rather
- than to address bias, since randomization already handles confounding.

#### 580 C.3 Observational Studies

- 581 If the data is not from an RCT, we move to observational methods. Here, we cannot rely on
- randomization to eliminate confounding. Techniques to address confounding varies according to the
- 583 characteristics of the data.

#### 584 C.3.1 Binary Treatments with Temporal Variation

- Difference in Differences (DiD) For binary treatments, we first consider Difference in Differences
- when time based confounding is a concern. DiD requires information about treatment timing, which
- could be either binary (pre or post treatment) or staggered (different units adopting treatment at
- 588 different time periods).
- A classic example is the minimum wage study by Card and Krueger [Card and Krueger, 1994].
- 590 Pennsylvania and New Jersey are two states with similar characteristics. New Jersey increased the
- minimum wage, but Pennsylvania did not. We want to know the effect of the minimum wage policy
- on employment. However, over time many things change that could affect employment outcomes.
- DiD uses the control group (Pennsylvania) as a counterfactual. Since the two states are similar, they
- would have evolved similarly absent treatment. By subtracting the changes in the control state from
- the changes in the treated state, DiD removes time varying confounders that affect both states equally.
- If no temporal information is available, DiD cannot be used.
- 597 Regression Discontinuity Design (RDD) If DiD is not applicable, we check for Regression
- 598 Discontinuity Design. RDD exploits situations where treatment assignment is determined by a
- threshold or cutoff rule, creating a sharp change in treatment status at a specific value of a running
- 600 variable.
- For example, to study the impact of alcohol consumption on road accidents [Carpenter and Dobkin,
- 602 2009], we can use the fact that alcohol consumption is legally allowed after age 21 in the United
- 603 States. By comparing accident rates for people just above and below age 21, we can identify the causal
- effect of legal drinking. Here, age is the running variable, and treatment assignment is: treatment = 1
- if age  $\geq 21$ , treatment = 0 if age < 21. If there is no clear running variable with a threshold that
- determines treatment assignment, RDD cannot be used.

# 607 C.3.2 General Methods for Binary Treatments

- 608 DiD and RDD can potentially work when the dataset meets certain characteristics. This would be
- the presence of treatment and outcome information over time for DiD and the presence of a running
- variable for RDD. If the dataset does not meet the respective criteria, we can rule those methods out.
- Next, we describe a more general group of methods.
- 612 **Backdoor adjustment** Observational studies often rely on the conditional ignorability assumption:
- treatment assignment is independent of potential outcomes conditional on all observed confounders.
- This requires that all relevant confounders are observed. We can satisfy this assumption using Pearl's
- backdoor adjustment criterion [Pearl, 2000]. For causal effect estimation, we use inverse probability
- 616 weighting (IPW) using propensity scores [Rosenbaum and Rubin, 1983] and matching [Stuart, 2010].

Propensity score measures the probability that a unit receives the treatment given the observed covariates. Mathematically, e(X) = P(T = 1|X). One can estimate propensity scores by fitting a logit or a probit model with X as the dependent variable and T as the independent variable. To choose between IPW and matching, we assess covariate balance using the standardized mean difference (SMD) [Greifer, 2025]:

$$\text{SMD} = \frac{\bar{X}_{\text{treated}} - \bar{X}_{\text{control}}}{\sqrt{\frac{s_{\text{treated}}^2 + s_{\text{control}}^2}{2}}}.$$

where  $\bar{X}_{\text{treated}}$  denotes the mean of the covariates in the treated group, and  $s_{\text{treated}}^2$  their variance. Analogously,  $\bar{X}_{\text{control}}$  and  $s_{\text{control}}^2$  refer to the mean and variance of the covariates in the control group.

If covariates are well balanced (typically SMD < 0.1), we use IPW methods directly. If covariates are poorly balanced, we use matching techniques to improve balance before estimating causal effects.

# 626 C.3.3 Methods for Unmeasured Confounding

Backdoor adjustment based methods fail if important confounders are unobserved. The presence of unobserved confounders decision is mostly made by combining domain knowledge and the data generating process. In cases where unobserved confounders are suspected, we can use Instrumental Variable or frontdoor estimation.

Instrumental Variables (IV) The first approach is IV analysis [Imbens, 2014], which requires a valid instrument satisfying two conditions:

- 1. **Relevance**: The instrument must be correlated with treatment ( $Cov(Z, T) \neq 0$ ).
- Exclusion restriction: The instrument should affect the outcome only through treatment, not directly.

The relevance assumption is testable, while the exclusion restriction must be justified through domain knowledge. A classic example of IV analysis is Card's study [Card, 1993] estimating returns to years of education, using geographic proximity to college as an instrument. Many unobserved factors (parental income, personal motivation) affect earnings and education, making it difficult to isolate education's causal effect. However, proximity to college affects educational attainment (relevance) but has no direct effect on earnings except through education (exclusion restriction). Finding good instruments is challenging, and quite often the dataset does not contain good candidates for an instrument.

Frontdoor Criterion Another option that may work in the presence of unobserved confounders is frontdoor estimation based on the frontdoor criterion [Pearl, 2000]. This works when there exists a mediator that completely captures the treatment's effect on the outcome and satisfies the frontdoor criterion:

- The mediator must completely intercept the path from the treatment to the outcome.
- The relationship between the treatment and the mediator must not be confounded.
- The treatment must block all confounding paths between the mediator and the outcome.

# 651 C.3.4 Method Hierarchy

633

634 635

636

637

638

639

640

642

643

648

649

650

When multiple methods are applicable, there is a clear preference hierarchy:

Instrumental Variables > Frontdoor Criterion > Backdoor Adjustment.

This hierarchy exists because IV and frontdoor methods can handle unobserved confounders, while backdoor adjustment requires that all confounders are observed. However, IV and frontdoor methods are less generally applicable, since finding valid instruments or suitable mediating variables is often challenging. Backdoor adjustment methods are more widely applicable but require stronger assumptions about confounding.

# C.3.5 Nonbinary treatments

For nonbinary treatments in observational studies, we consider three methods: instrumental variables (IV), frontdoor adjustment, and generalized propensity scores using the backdoor adjustment set. IV and frontdoor approaches apply almost exactly as in the case of binary treatments described above. Backdoor estimation also works similarly, except that we use a different estimation method to compute causal effects. In this case, we use generalized propensity scores, which extend the idea of propensity scores to continuous treatments [Hirano and Imbens, 2004]. The scores are computed using the backdoor adjustment set.

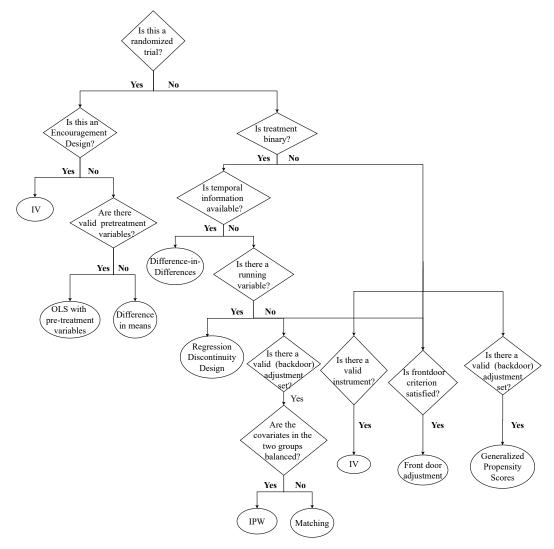


Figure 3: Decision-tree that guides method selection in CAIS. We prompt an LLM to generate responses to queries corresponding to the decision nodes, and traverse the tree accordingly before reaching a leaf node, which corresponds to a method.

# 566 D Detailed Methodology

In this section we extend the Methodology of CAIS discussed in Section 4 in detail. CAIS operates in four sequential stages. Each stage comprises one or more micro-tools that pass a typed artifact to the next stage. A validation loop connects Stage 3 back to Stage 2 when assumptions fail. In total, CAIS uses eight micro-tools: (1) input\_parser, (2) dataset\_analyzer, (3) query\_interpreter, (4) method\_selector, (5) method\_validator, (6) method\_executor, (7) explanation\_generator, and (8) output\_formatter. All tools read/write a shared typed state (dataset profile, normalized query, analysis plan, selected method with assumptions, diagnostics, estimates, and final report).

#### 674 D.1 Stage 1: Data Preprocessing & Query Decomposition

We begin by using **input\_parser** tool to analyze user-specified causal queries and extract three 675 key components: the *query type*, the *relevant variables*, and any *explicit constraints*. It adopts a 676 677 hybrid strategy combining regex-based heuristics with LLM-driven structured parsing. A tailored prompt E.1, enriched with dataset context when available, guides the LLM to classify the query 678 (e.g., EFFECT\_ESTIMATION, COUNTERFACTUAL, CORRELATION, DESCRIPTIVE, or OTHER) and to output 679 a structured JSON conforming to a predefined schema. This schema enforces explicit roles for 680 variables (treatment, outcome, covariates, grouping variables, instruments) and records constraints or 681 dataset paths if mentioned. The LLM output is validated to ensure logical consistency (e.g., requiring both treatment and outcome for effect estimation queries) and alignment with dataset columns when 683 available. Regex-based matching serves as a complementary mechanism, particularly for dataset path extraction, thereby improving robustness. The final output is a standardized dictionary encapsulating 685 the original query, its type, extracted variables, constraints, and dataset path. 686

Following this, the **dataset\_analyzer** tool profiles the dataset to identify characteristics relevant 687 for causal inference. It begins by extracting basic metadata such as the number of rows, columns, 688 file name, and data types. A detailed analysis is then conducted to detect temporal structures, 689 panel data patterns, and potential discontinuities, alongside computing correlations among numeric 690 variables. The module further identifies *potential treatments* and *outcomes*, either heuristically or 691 with LLM assistance, and assesses candidate instrumental variables. The prompts used in this stage 692 are relatively light weight and mainly framed as direct questions to the LLM, an example prompt 693 for identifying potential treatment and outcome variables can be found in Appendix E.2. For binary 694 treatment candidates, it computes per-group summary statistics (e.g., group sizes, means, and standard 695 deviations of covariates) to facilitate balance checks. Additionally, it records missing values, unique 696 value counts, and column categorizations to provide a structured overview of the dataset. The final 697 output is a comprehensive dictionary that consolidates dataset information, candidate causal variables, detected structures, and diagnostic statistics. 699

Finally, the query\_interpreter reconciles the normalized query with the dataset profile to materi-700 alize the treatment variable T, outcome Y, admissible covariates X, and any design-specific fields 701 (instruments Z, running variable R with cutoff c, time and group indicators for DiD, etc.). For each of 702 these fields, the system issues targeted prompts to an LLM, these structured prompts can be found in 703 Appendix E.3. The responses are aggregated and validated to ensure coherence with both the dataset 704 schema and causal design assumptions. The final output is a structured analysis plan comprising 705 named columns, encoding details (e.g., reference levels, interaction terms), and a canonical estimand (such as ATE, ATT, LATE, sharp/fuzzy RDD, or DiD horizon), which then serves as the input for 707 subsequent method selection phase. 708

# D.2 Stage 2: Method Selection

709

Given the analysis plan, the **method\_selector** chooses a candidate estimator via a deterministic decision tree that encodes standard causal design logic. The selector also emits an assumption checklist and an ordered fallback list to support backtracking. This stage is completely deterministic and requires no involvement of LLMs. The design choice behind the nodes of the decision tree is already discussed in great detail in Appendix C.

#### 715 D.3 Stage 3: Validation

726

Once a candidate estimator is selected, the **method\_validator** module verifies whether the assump-716 tions required for its validity hold in the given dataset. This includes statistical diagnostics (e.g., 717 overlap and positivity checks for treatment assignment, relevance and exogeneity tests for instruments, 718 bandwidth support for RDD, and parallel trends diagnostics for DiD). The validator also cross-checks 719 variable roles (treatment, outcome, covariates, instruments, etc.) against the dataset schema to ensure 720 consistency. Failures in any diagnostic trigger a backtracking mechanism: the system reverts to 721 the Method Selector's fallback list to propose an alternative estimator whose assumptions are more 722 compatible with the data. This creates a validation loop that guarantees every executed method 723 is grounded in both causal theory and empirical feasibility, while preserving determinism in the 724 diagnostic procedures. In this stage as well no LLMs are involved.

#### D.4 Stage 4: Method Execution & Interpretation

On a validated design, the Method Executor runs a pre-verified estimator template (statsmodels/2SLS, 727 DiD with appropriate fixed effects, local polynomial RDD, matching/weighting for PSM/PSW, 728 GPS, diff-in-means). Execution favors templates over code-generation for stability and speed. One cross-cutting prompt, STATSMODELS\_PARAMS\_IDENTIFICATION\_PROMPT\_TEMPLATE, helps select the correct coefficient(s) to report when formulas include encodings or interactions (used by linear/GLM-style estimators). Several methods expose small, function-local LLM assists 732 for parameter suggestions or narrative-e.g., IV, RDD, DiD, GPS/backdoor/diff-in-means/linear 733 regression helpers, each constrained to strict JSON and backed by deterministic fallbacks. Finally 734 Explanation Generator then converts structured artifacts (chosen design, diagnostics, estimates) into 735 a concise, dataset-specific justification and interpretation, while the Output Formatter assembles the 736 final response object with numeric results, uncertainty, method card (assumptions & checks), and a 737 short plain-language answer for downstream rendering. 738

# 739 E Methodology Prompts

## 40 E.1 Causal Query Parsing Prompt

Analyze the following causal query **strictly in the context of the provided dataset information (if available)**. Identify the query type, key variables (mapping query terms to actual column names when possible), constraints, and any explicitly mentioned dataset path. User Query: "{query}

{dataset\_context}

# **Guidance for Identifying Query Type:**

- **EFFECT\_ESTIMATION**: Look for keywords like "effect", "impact", "influence", "cause", "affect", "consequence". Also consider questions asking "how does X affect Y?" or comparing outcomes between groups based on an intervention.
- **COUNTERFACTUAL**: Look for hypothetical scenarios, often using phrases like "what if", "if X had been", "would Y have changed", "imagine if", "counterfactual".
- **CORRELATION**: Look for keywords like "correlation", "association", "relationship", "linked to", "related to". These queries ask about statistical relationships without necessarily implying causality.
- **DESCRIPTIVE**: Queries that ask for summaries, descriptions, trends, or statistics about the data without investigating causal links (e.g., "Show sales over time", "What is the average age?").
- OTHER: Use if the query does not fit any of the above categories.

Choose the most appropriate type from: EFFECT\_ESTIMATION, COUNTERFACTUAL, CORRELATION, DESCRIPTIVE, OTHER.

#### Variable Roles to Identify:

• treatment: The intervention or variable whose effect is being studied.

- outcome: The result or variable being measured.
- covariates\_mentioned: Variables explicitly mentioned to control for or adjust for.
- grouping\_vars: Variables identifying specific subgroups for analysis (e.g., "for men", "in the sales department").
- instruments mentioned: Variables explicitly mentioned as potential instruments.

**Constraints:** Conditions applied to the analysis (e.g., filters on columns, specific time periods).

**Dataset Path Mentioned:** Extract the file path or URL if explicitly stated in the query.

**Output:** ONLY a valid JSON object matching this schema (no explanations or surrounding text):

```
"query_type": "<Identified Query Type>",
"variables": {
    "treatment": ["<Treatment Variable(s) Mentioned>"],
        "outcome": ["<Outcome Variable(s) Mentioned>"],
        "covariates_mentioned": ["<Covariate(s) Mentioned>"],
        "grouping_vars": ["<Grouping Variable(s) Mentioned>"],
        "instruments_mentioned": ["<Instrument(s) Mentioned>"]
},
"constraints": ["<Constraint 1>", "<Constraint 2>"],
        "dataset_path_mentioned": "<Path Mentioned or null>"
```

If Dataset Context is provided, ensure variable names in the output JSON correspond to actual column names where possible. If no context is provided, or if a mentioned variable doesn't map directly, use the phrasing from the query. Respond with only the JSON object.

742

# **E.2** Light Weight Inline Prompts for Dataset Analysis

#### **Prompt:**

You are an expert causal inference data scientist. Identify potential treatment and outcome variables from this dataset.

```
Dataset Description: {description_text}
Dataset Columns: {columns_list}
Column Types: {column_types}
```

**Binary Columns (good treatment candidates):** {binary\_cols}

#### **Instructions:**

- 1. Identify **TREATMENT** variables: interventions, treatments, programs, policies, or binary state changes. Look for binary variables or names with "treatment", "intervention", "program", "policy", etc.
- 2. Identify **OUTCOME** variables: results, effects, or responses to treatments. Look for numeric variables (especially non-binary) or names with "outcome", "result", "effect", "score", etc.

```
Output: ONLY a valid JSON object with two lists:
```

```
{
   "potential_treatments": ["treatment_a", "program_b"],
   "potential_outcomes": ["result_score", "outcome_measure"]
}
```

# **Instrumental Variable Identification Prompt**

#### **Prompt:**

You are a causal inference assistant tasked with assessing whether a valid Instrumental Variable (IV) exists in the dataset. A valid IV must satisfy **all** of the following conditions:

- 1. **Relevance**: It must causally influence the Treatment.
- Exclusion Restriction: It must affect the Outcome only through the Treatment —
  not directly or indirectly via other paths.
- 3. **Independence**: It must be as good as randomly assigned with respect to any unobserved confounders affecting the Outcome.
- 4. **Compliance (for RCTs)**: If the dataset comes from a randomized controlled trial or experiment, IVs are only valid if compliance data is available i.e., if some units did not follow their assigned treatment. In this case, the random assignment may be a valid IV, and compliance is the actual treatment variable. If compliance related variable is not available, do not select IV.
- 5. The instrument must be one of the listed dataset columns (not the treatment itself), and must not be assumed or invented.

You should **only suggest an IV if you are confident that all the conditions are satisfied**. Otherwise, return "NULL".

#### **Information Provided:**

• User Query: {query}

• Dataset Description: {description}

• Treatment: {treatment}

• Outcome: {outcome}

• Available Columns: {column\_info}

**Output:** Return a JSON object with the following structure:

{ "instrument\_variable": "COLUMN\_NAME\_OR\_NULL" }

747

# **RDD Identification Prompt**

#### **Prompt:**

You are an expert causal inference assistant helping to determine if Regression Discontinuity Design (RDD) is applicable for quasi-experimental analysis.

#### **Information Provided:**

- User Query: {query}
- Dataset Description: {description}
- Identified Treatment (tentative): {treatment}
- Identified Outcome (tentative): {outcome}
- Available Columns: {column\_info}

**Task:** Your goal is to check if there is a **Running Variable**, i.e., a variable that determines treatment/control.

- If the variable is above a certain cutoff, the unit is categorized as treated; if below, it is control.
- The running variable must be numeric and continuous. Do not use categorical or low-cardinality variables.

• The treatment variable must be binary. If not, RDD is not valid.

**Output:** Respond ONLY with a valid JSON object. If RDD is not suggested by the context, return null for both fields.

```
{ "running_variable": "COLUMN_NAME_OR_NULL", "cutoff_value": NUMERIC_VALUE_OR_NULL }

Examples:
{ "running_variable": "test_score", "cutoff_value": 70 }
{ "running_variable": null, "cutoff_value": null }
```

749

# **RCT Identification Prompt**

#### **Prompt:**

You are an expert causal inference assistant helping to determine if the data comes from a Randomized Controlled Trial (RCT). Your goal is to assess if the treatment assignment mechanism described or implied was random.

#### **Information Provided:**

- User Query: {query}
- Dataset Description: {description}
- Identified Treatment (tentative): {treatment}
- Identified Outcome (tentative): {outcome}
- Available Columns: {column\_info}

**Output:** Respond ONLY with a valid JSON object matching the required schema.

```
{ "is_rct": BOOLEAN_OR_NULL }
Examples:
{ "is_rct": true } # RCT likely
{ "is_rct": false } # Observational likely
{ "is_rct": null } # Unsure
```

750

# **Treatment Reference Identification Prompt**

#### Promnt

You are a causal inference assistant.

#### **Dataset Information:**

- Dataset Description: {description}
- Identified Treatment Variable: "{treatment\_variable}"
- Unique Values in Treatment Variable (sample): {treatment\_variable\_values}
- User Query: {query}

**Task:** Based on the user query, determine if it specifies a particular category of the treatment variable {treatment\_variable} that should be considered the control, baseline, or reference group for comparison.

# **Examples:**

- Query: "Effect of DrugA vs Placebo" → Reference for treatment "Drug" = "Placebo"
- Query: "Compare ActiveLearning and StandardMethod against NoIntervention" → Reference for "TeachingMethod" = "NoIntervention"

If a reference level is clearly specified or strongly implied **and** it is one of the unique values provided, identify it. Otherwise, state null. If multiple values seem like controls (e.g., "compare A and B vs C and D"), return null for now.

```
Output: Return ONLY a JSON object adhering to this schema:
{
    "reference_level": "string_representing_the_level_or_null",
    "reasoning": "string_or_null_brief_explanation"
}
```

752

# **Interaction Term Identification Prompt**

#### **Prompt:**

You are a causal inference assistant.

Your task is to determine whether the user query suggests the inclusion of an interaction term between the treatment and one covariate, specifically to assess heterogeneous treatment effects (HTE).

# **Information Provided:**

- User Query: {query}
- Dataset Description: {description}
- Identified Treatment Variable: {treatment\_variable}
- Available Covariates (name: type): {covariates\_list\_with\_types}

#### **Instructions:**

- ONLY suggest an interaction if the query explicitly mentions treatment across a subgroup.
- DO NOT suggest an interaction if the query asks for an overall average effect or does not mention subgroup analysis.
- If unsure, default to no interaction.

```
Output Schema:
```

# **Treatment Variable Identification Prompt**

## **Prompt:**

You are an expert in causal inference. Your task is to identify the **treatment variable** in a dataset in order to perform a causal analysis that answers the user's query.

#### **Information Provided:**

• User Query: {query}

• Dataset Description: {description}

• List of Available Variables: {column\_info}

**Task:** Based on the query, dataset description, and available variables, determine which variable is most likely to serve as the treatment variable.

If a clear treatment variable cannot be determined, return null.

## **Output Schema:**

```
{ "treatment": "COLUMN_NAME_OR_NULL" }
```

754

# **Outcome Variable Identification Prompt**

#### **Prompt:**

You are an expert in causal inference. Your task is to identify the **outcome variable** in a dataset in order to perform a causal analysis that answers the user's query.

#### **Information Provided:**

• User Query: {query}

• Dataset Description: {description}

• Available Variables: {column\_info}

**Task:** Based on the query, dataset description, and available variables, determine which variable is most likely to serve as the outcome variable in the causal analysis. Do not speculate. If a clear outcome variable cannot be identified, return null.

#### **Output Schema:**

```
{ "outcome": "COLUMN_NAME_OR_NULL" }
```

755

# **Covariates Identification Prompt**

#### **Prompt:**

You are an expert in causal inference. Your task is to identify the **pre-treatment variables** in a dataset that can be used as controls in a causal estimation model to answer the user's query.

#### **Information Provided:**

• User Query: {query}

• Dataset Description: {description}

• Available Variables: {column\_info}

• Treatment Variable: {treatment}

• Outcome Variable: {outcome}

**Task:** Pre-treatment variables are those that are measured **before** the treatment is applied and are **not affected** by the treatment. These variables can be used as controls in the causal model.

For example, in an RCT with outcome Y, treatment T, and pre-treatment variables X1, X2, X3, we can perform a regression of the form:

$$Y \sim T + X1 + X2 + X3$$

Based on the information above, return a list of variables that qualify as pre-treatment variables from the available columns. If no suitable pre-treatment variables can be identified, return an empty list.

# **Output Schema:**

```
{ "covariates": ["LIST_OF_COLUMN_NAMES_OR_EMPTY_LIST"] }
```

757

# **Estimand Identification Prompt**

## **Prompt**:

You are an expert in causal inference. Your task is to determine the appropriate **estimand** to answer a given query.

## **Information Provided:**

• User Query: {query}

• Dataset Description: {dataset\_description}

• Variables in Dataset: {dataset\_columns}

• Treatment Variable: {treatment}

• Outcome Variable: {outcome}

**Task:** Given this information, decide whether the **Average Treatment Effect (ATE)** or the **Average Treatment Effect on the Treated (ATT)** is more appropriate for answering the query.

**Output:** Only return the estimand name:

"att" or "ate"

758

# **Confounder Identification Prompt**

#### **Prompt**:

You are an expert in causal inference. Your task is to identify potential **confounders** in a dataset that should be adjusted for when estimating the causal effect described in the user query.

#### **Information Provided:**

• User Query: {query}

• Dataset Description: {description}

• Available Variables: {column\_info}

• Treatment Variable: {treatment}

• Outcome Variable: {outcome}

#### **Definition of Confounder:** A **confounder** is a variable that:

- 1. Affects the treatment (influences who receives the treatment),
- 2. Affects the outcome,
- 3. Is not caused by the treatment (must be a pre-treatment variable),
- 4. Is not a mediator between treatment and outcome.

These variables can create spurious associations between treatment and outcome if not adjusted for.

**Task:** Based on the user query and the dataset description, identify which variables are likely to be confounders. Only include variables that you believe causally affect both treatment and outcome. If uncertain, only include variables where the justification is clear from the query or description.

# **Output Schema:**

```
{ "confounders": ["LIST_OF_COLUMN_NAMES_OR_EMPTY_LIST"] }
```

760

# **DiD Term Identification Prompt**

#### **Prompt:**

You are a causal inference assistant tasked with determining whether a valid Difference-in-Differences (DiD) **interaction term** already exists in the dataset.

This DiD term should be a **binary variable** indicating whether a unit belongs to the **treatment group after treatment was applied**.

For example, if a policy was enacted in 2020 for a particular state, then the DiD term would equal 1 for units from that state in years after 2020, and 0 otherwise.

#### **Information Provided:**

• User Query: {query}

• Time Variable: {time\_variable}

• Group Variable: {group\_variable}

• Dataset Description: {description}

• Available Columns: {column\_info}

• Column Types: {column\_types}

**Output:** Return your answer as a valid JSON object with the following format:

```
{ "did_term": "COLUMN_NAME_OR_NULL" }
```

761

#### 2 F Detailed Study: Method Validation Loop

This example presents the complete prompt employed in the validation feedback loop of CAIS.

# **Worked Example: Method Validation**

**Query:** Does having access to electricity increase kerosene expenditures?

Dataset: electrification\_data.csv

**Database:** All\_Data Collection (Rural Electrification Survey)

**Description:** This household survey covers 686 households in 120 habitations across Uttar Pradesh, India. Using a geographic eligibility rule (households within 20–35 m vs. 45–60 m of a power pole), it records monthly expenditures on food, education, kerosene, total expenditure, appliance ownership, lighting usage, and satisfaction measures to assess the impact of electrification.

**Method Validation:** During validation, the pipeline fits local regressions on kerosene expenditure immediately below and above the 40 m cutoff to test for a sharp discontinuity. When using the lightweight gpt-4o-mini model, the agent misidentified the "distance" variable effectively widening the window around 40 m and consequently observed no statistically

significant jump in outcomes at the threshold (p>0.05). Because a pronounced, localized shift at the cutoff is the cornerstone of RDD, this absence of any detectable discontinuity constituted a direct violation of the RDD assumptions and led to its rejection. The system then automatically backtracked down the decision tree, removed RDD from consideration, and evaluated the next class of methods. Given the observational nature of the data and the rich set of covariates, it advanced to propensity-score-matching as the alternative method to create balanced treatment and control groups before estimating the effect.

765

766

# **G** Baselines Prompts

Here, we display all the prompts used for the baselines: ReAct, PoT, and Veridical Prompts for causal inference.

# **ReAct Prompt Example**

**Prompt:** You are working with a pandas DataFrame in Python. The name of the DataFrame is df.

You should use the tools below to answer the question posed to you:

python\_repl\_ast: A Python shell. Use this to execute Python commands. Input should be a valid Python command. When using this tool, sometimes output is abbreviated—make sure it does not look abbreviated before using it in your answer.

## Use the following format:

- Question: the input question you must answer
- Thought: what you should do next
- Action: the action to take (e.g., python\_repl\_ast)
- Action Input: the input to the action (code to execute)
- Observation: the result of the action

(This Thought/Action/Action Input/Observation can repeat N times.)

**Final Answer:** The final answer to the original input question. Please provide a structured response including the following:

- Method
- · Causal Effect
- Standard Deviation
- Treatment Variable
- Outcome Variable
- Covariates
- Instrument / Running Variable / Temporal Variable
- · Results of Statistical Test
- Explanation for Model Choice
- Regression Equation

#### **Instructions:**

- Import libraries as needed.
- Do not create any plots.
- Use the print() function for all code outputs.
- If you output an Action step, stop after generating the Action Input and await execution.
- If you output the Final Answer, do not include an Action step.

# Example Usage of python\_repl\_ast:

Action: python\_repl\_ast

770

# **Program of Thoughts based Prompt**

**Prompt:** You are a data analyst with strong quantitative reasoning skills. Your task is to answer a data-driven causal question using the provided dataset. The dataset description and query are given below.

You should analyze the **first 10 rows** of the dataset and then write Python code to generalize the analysis to the full table. You may use any Python libraries.

The returned value of the program should be the final answer. Please follow this format:

```
def solution():
    # import libraries if needed
    # load data from {self.dataset_path}
    # write code to get the answer
    # return answer
print(solution())
```

Dataset Description: {self.dataset\_description} Dataset Path:
{self.dataset\_path}

First 10 rows of data: {df.head(10)}

**Question:** {self.query}

# **Example Methods (choose one if applicable):**

- propensity\_score\_weighting: output the ATE
- propensity\_score\_matching\_treatment\_to\_control: output the ATT
- linear\_regression: output coefficient of variable of interest
- instrumental\_variable: output coefficient
- matching: output the ATE
- difference\_in\_differences: output coefficient
- regression\_discontinuity\_design: output coefficient
- linear\_regression / difference\_in\_means: output coefficient / DiM

**Response:** The final answer should include a structured summary with the following fields (use "NA" where not applicable):

- Method
- · Causal Effect
- · Standard Deviation
- Treatment Variable
- Outcome Variable
- Covariates
- Instrument / Running Variable / Temporal Variable
- · Results of Statistical Test
- Explanation for Model Choice
- Regression Equation

## **Veridical Prompt**

You are an expert in statistics and causal reasoning. You will use a rigorous scientific framework to answer a causal question using a structured, step-by-step process with checklists. Problem Statement: self.query

**Step 1: Domain Understanding** - What is the real-world question? Why is it important? - Could alternate formulations impact the final result?

**Step 2: Dataset Overview** - Dataset Path: dataset\_path - Description: dataset\_description - Dataset Summary, Types, Missing Values, Preview Rows

Checklist: - How was data collected? Design principles? - What are the variables, types, and units? - Are there errors or pre-processing artifacts?

**Step 3: Exploratory Analysis** - Identify confounders, mediators, biases - Suspect endogeneity? What instruments might be relevant? - Are strong correlations present?

**Step 4: Modeling Strategy** - Choose 3 candidate methods (e.g., matching, regression, IV) - State assumptions and reasons for each method - Discuss software libraries to be used and potential pitfalls - Outline key outputs and steps in analysis

**Step 5: Post Hoc Analysis** - Are relationships or outcomes unexpected? - Assess result stability and robustness

**Step 6: Interpretation and Reporting** Final Answer: Report the following fields: - Method, Causal Effect, Standard Deviation - Treatment and Outcome Variables - Covariates, Instruments or Temporal Elements - Results of any statistical tests - Justification of model choice - Equation or summary used