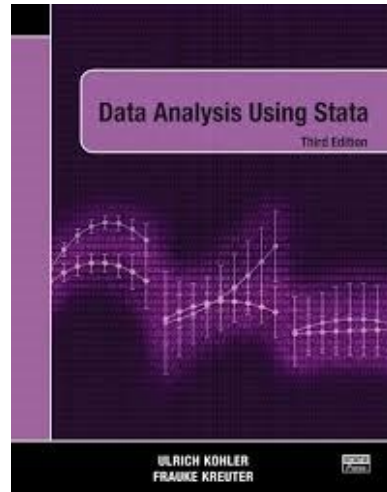
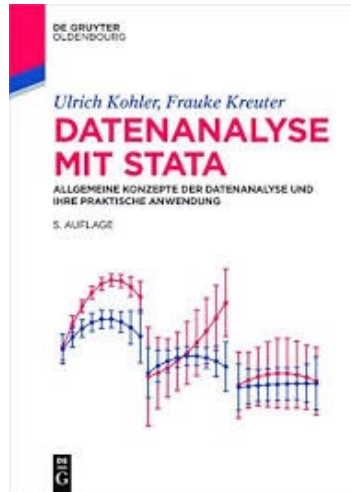


Cocreating with AI: LLMs as data scientist agents



Hamburg May 28 2025



Data Analysis
Using AI

4.4.4 LLMs as programming assistant for Python

As outlined in Chapter [1.3.15](#), Large Language Models (LLMs) are valuable tools for coding assistance, particularly with languages like Python, where they have shown considerable efficacy. Suppose we are in the Python environment within Stata, having imported the variable 'ybirth' using the `sfi` module. We seek to perform further computations on this data and are turning to Python programming for assistance. We cannot assure that the prompts and requests demonstrated here will be effective with every LLM or under all circumstances, but we will point out several key features of the prompt that have helped us in the past to get better results:

I am working with the already existing variable ybirth and I want to calculate its mean in Python without any additional modules. The code should all be put into one code section, simple, well-commented, and easy to understand, and reusable. It is important that it is efficient and follows best practices. The goal is to create a solid foundation for analysis that can be easily understood and adapted by others.

Perspective on Data Science

Roger D. Peng¹ and Hilary S. Parker²

¹Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland 21205, USA; email: rdpeng@jhu.edu

²Independent Consultant, San Francisco, California 94102, USA

Keywords

data analysis, analytic iteration, design thinking, reproducibility, systems engineering

Abstract

The field of data science currently enjoys a broad definition that includes a wide array of activities which borrow from many other established fields of study. Having such a vague characterization of a field in the early stages might be natural, but over time maintaining such a broad definition becomes unwieldy and impedes progress. In particular, the teaching of data science is hampered by the seeming need to cover many different points of interest. Data scientists must ultimately identify the core of the field by determining what makes the field unique and what it means to develop new knowledge in data science. In this review we attempt to distill some core ideas from data science by focusing on the iterative process of data analysis and develop some generalizations from past experience. Generalizations of this nature could form the basis of a theory of data science and would serve to unify and scale the teaching of data science to large audiences.

Annu. Rev. Stat. Appl. 2022. 9:1–20

First published as a Review in Advance on September 16, 2021

The *Annual Review of Statistics and Its Application* is online at [statistics.annualreviews.org](https://doi.org/10.1146/annurev-statistics-040220-013917)

<https://doi.org/10.1146/annurev-statistics-040220-013917>

Copyright © 2022 by Annual Reviews.
All rights reserved

ANNUAL
REVIEWS **CONNECT**

www.annualreviews.org

- Download figures
- Navigate cited references
- Keyword search
- Explore related articles
- Share via email or social media

arXiv:2408.07461v1 [cs.AI] 14 Aug 2024

Problem Solving Through Human-AI Preference-Based Co-operation

Subhabrata Dutta
TU Darmstadt

subhabrata.dutta@tu-darmstadt.de

Timo Kaufman
LMU Munich

timo.kaufmann@lmu.de

Goran Glavas
University of Würzburg

goran.glavas@uni-wuerzburg.de

Ivan Habernal
RU Bochum

ivan.habernal@ruhr-uni-bochum.de

Kristian Kersting
TU Darmstadt

kersting@cs.tu-darmstadt.de

Frauke Kreuter
LMU Munich

frauke.kreuter@lmu.de

Mira Mezini
TU Darmstadt

mezini@informatik.tu-darmstadt.de

Iryna Gurevych
TU Darmstadt

iryna.gurevych@tu-darmstadt.de

Eyke Hüllermeier*
LMU Munich

eyke@ifi.lmu.de

Hinrich Schuetze*
LMU Munich

hhschutze@gmx.de

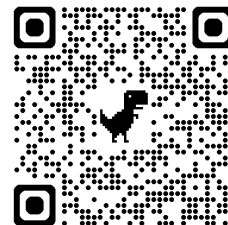
Abstract

While there is a widespread belief that artificial general intelligence (AGI) – or even super-human AI – is imminent, complex problems in expert domains are far from being solved. We argue that such problems require human-AI cooperation and that the current state of the art in generative AI is unable to play the role of a reliable partner due to a multitude of shortcomings, including inability to keep track of a complex solution artifact (e.g., a software program), limited support for versatile human preference expression and lack of adapting to human preference in an interactive setting. To address these challenges, we propose HAI-Co², a novel human-AI co-construction framework. We formalize HAI-Co² and discuss the difficult open research problems that it faces. Finally, we present a case study of HAI-Co² and demonstrate its efficacy compared to monolithic generative AI models.

1 Introduction

Despite the impressive achievements of generative AI spearheaded by Large Language Models (LLMs), Language Models and code models (Lozhkov et al., 2024; Wang et al., 2021), multiple recent investi

*Shared senior authorship



Tukey defined data analysis as
“procedures for **analyzing data**,
techniques for **interpreting the results** of such procedures,
ways of planning the **gathering of data** to make its analysis
easier, more precise or more accurate, and all the
machinery and results of (mathematical) statistics which
apply to analyzing data” ([**Tukey 1962**](#), p. 2)

The Excitement

Tukey defined data analysis as
“procedures for analyzing data,
techniques for interpreting the results of such procedures,
ways of planning the gathering of data to make its analysis
easier, more precise or more accurate, and all the
machinery and results of (mathematical) statistics which
apply to analyzing data” (Tukey 1962, p. 2)

AI analyzing data



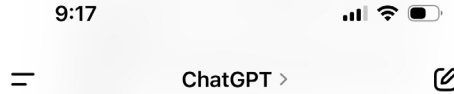
I want to use the Titanic dataset to explore with my students that question whether females and children were prioritized in being saved in the Titanic accident. Can you load the data and show me how to do this in Stata?



Yes! The Titanic dataset is a great resource for exploring whether women and children were prioritized in survival. I'll guide you through loading the dataset and conducting the analysis in Stata.

Step 1: Load the Dataset in Stata

You can download the Titanic dataset from sources like Kaggle or use Stata's built-in datasets (if available). The dataset typically includes variables such as:



- sex: Male or Female
- age: Age of the passenger

If you have a CSV file (e.g., titanic.csv), load it into Stata using:

```
import delimited "titanic.csv", clear
```

Step 2: Check the Data

Verify that your dataset loaded correctly:

```
describe  
list survived sex age if _n <= 10 //  
Preview first 10 rows
```

If age has missing values, check:

```
sum age, detail
```

Step 3: Analyze Survival Rates

To explore whether women and children were prioritized, we can compare

2. Compare Survival by Age (Defining "Children")

First, create a variable for children (e.g., age < 18):

```
gen child = age < 18
```

Then check survival by gender and child status:

```
tab child survived, col  
tab sex child survived, col
```

3. Logistic Regression: Were Women and Children More Likely to Survive?

A logistic regression can help quantify the likelihood of survival:

4. Interaction Between Gender and Being a Child

To check if the effect of being a child differs by gender:

```
logit survived i.sex##i.child i.pclass
```

Step 4: Interpret Results

- If the coefficient for sex is positive and significant, women had higher odds of survival.
- If child has a positive coefficient, children were prioritized.

AI analyzing data



Good question! Here are some refinements that could improve your analysis:

1. Adjust for Socioeconomic Status (Passenger Class)

- We know that **passenger class (pclass)** played a major role in survival.
- Higher-class passengers had better access to lifeboats.

Refinement: Add pclass as a Control Variable

Modify your logistic regression to adjust for class differences:

```
logit survived i.sex##i.child i.pclass
```

- This helps separate the effects of



Refinement: Visualize Survival by Class

```
graph bar (mean) survived, over(sex)  
over(child) over(pclass), blabel(bar)
```

- This helps you see if class influenced survival within each group.

2. Consider Non-Linearity in Age Effects

Instead of a simple **child vs. adult** binary, age could have a more **gradual effect**.

Refinement: Use Age as a Continuous Variable

```
logit survived i.sex c.age i.pclass
```

- Here, `c. age` treats age as a continuous variable.
- Check whether the effect of age on survival is linear

Refinement: Check for a Nonlinear Effect (Quadratic Term)

```
gen age2 = age^2  
logit survived i.sex c.age c.age2  
i.pclass
```

- If `age2` is significant, age has a **curved** effect on survival.

3. Interaction Between Age and Class

- Were **children in lower classes** at a disadvantage?
- Test if the effect of being a child varies by class.

Refinement: Interaction Between Child and Class

```
logit survived  
i.sex##i.child##i.pclass
```

- This tests if **children in 3rd class had lower survival than children in 1st class**.

Refinement: Marginal Effects by Class

```
margins child#pclass  
marginsplot, xlabel(1 "1st Class" 2  
"2nd Class" 3 "3rd Class")
```

- This shows survival rates for **children vs. adults across classes**.

Tukey defined data analysis as
“procedures for analyzing data,
techniques for interpreting the results of such procedures,
ways of planning the **gathering of data** to make its analysis
easier, more precise or more accurate, and all the
machinery and results of (mathematical) statistics which
apply to analyzing data” (Tukey 1962, p. 2)



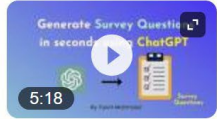
Showing results for designing questionnaires with **chatgpt**
Search instead for designing questionnaires with chaptgpt

📺 Videos :



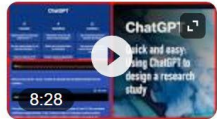
How to design a survey using ChatGPT

YouTube · Mohamed Benhima
Apr 16, 2023



How to generate survey questions for research using ChatGPT

YouTube · Tasvir Mahmood
Mar 26, 2023



Using ChatGPT to design a research study: research question ...

YouTube · Ranywayz Random
Jan 28, 2023



10 key moments in this video ▼



How to use ChatGPT to create survey items: 3 examples ...

YouTube · Ranywayz Random
Feb 11, 2023



3 key moments in this video ▼

1. **On a scale of 1-10, how satisfied were you with this class?**
2. **How relevant and useful did you find the content of the class for your academic and/or professional goals?**
(1 = Not relevant at all, 10 = Extremely relevant)
3. **How effective was the instructor in delivering the course material?**
(1 = Not effective at all, 10 = Extremely effective)
4. **How encouraged did you feel to participate and engage in class discussions?**
(1 = Not encouraged at all, 10 = Highly encouraged)

Improve this question “On a scale of 1-10, how satisfied were you with this class?”
Follow the textbook from Schuman and Presser.
Structure output “Revised Question:” “Answer Options:” and Explanation



Revised Question: How would you describe your overall satisfaction with this class?

Answer Options: Very dissatisfied, Dissatisfied, Neutral, Satisfied, Very satisfied.

Explanation: A 5-point scale simplifies the choices and captures nuanced feelings.

Avoiding numbers ensures that participants don't just choose a 'middle' number without giving it thought.

And this one: “How relevant and useful did you find the content of the class for your academic and/or professional goals? (1 = Not relevant at all, 10 = Extremely relevant)”



How would you rate the relevance of the class content to your academic or professional goals? **Answer Options:** Not at all relevant, Slightly relevant, Moderately relevant, Very relevant, Extremely relevant. *Explanation: Splitting relevance and usefulness makes it clearer and avoids confusion. Simplified options provide clarity.*



[Political Analysis](#)

Article contents

[Abstract](#)

[Footnotes](#)

[References](#)

Out of One, Many: Using Language Models to Simulate Human Samples

Published online by Cambridge University Press: 21 February 2023

Lisa P. Argyle , Ethan C. Busby, Nancy Fulda, Joshua R. Gubler , Christopher Rytting and David Wingate

[Show author details](#) ▼

Article


[Supplementary materials](#)

[Metrics](#)

[Get access](#)

 [Share](#)

 [Cite](#)

 [Rights & Permissions](#)

Abstract

We propose and explore the possibility that language models can be studied as effective proxies for specific human subpopulations in social science research. Practical and research applications of artificial intelligence tools have sometimes been limited by problematic biases (such as racism or sexism), which are often treated as uniform properties of the models. We show that the “algorithmic bias” within one such tool—the GPT-3 language model—is instead both fine-grained and demographically correlated, meaning that proper conditioning will cause it to accurately emulate response distributions from a wide variety of human subgroups. We term this property *algorithmic fidelity* and explore its extent in GPT-3. We create “silicon samples” by conditioning the model on thousands of sociodemographic backstories from real human participants in multiple large surveys conducted in the United States. We then compare the silicon and human samples to demonstrate that the

13

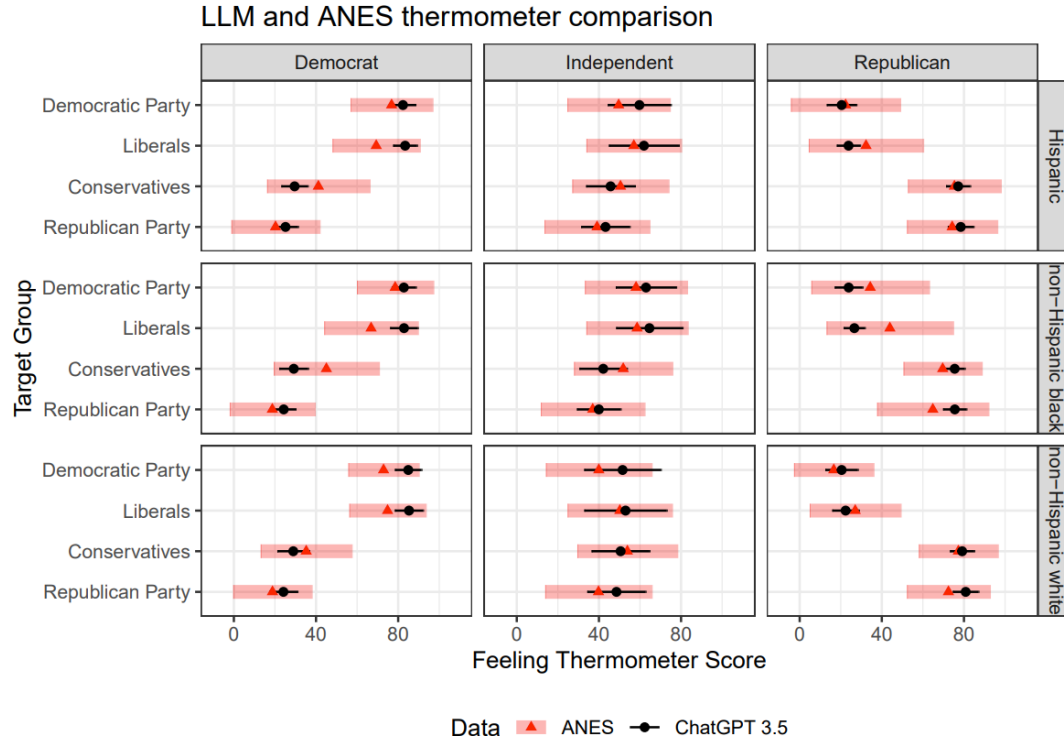


Figure 2: Average feeling thermometer results (x-axis) for different target groups (y-axes) by party ID of respondent (columns). Average ANES estimates from the 2016 and 2020 waves indicated with red triangles and one standard deviation indicated with thick red bars. LLM-derived averages indicated by black circles and thin black bars. Sample sizes for each group-wise comparison are identical.

Bisbee, J., Clinton, J., Dorff, C., Kenkel, B., & Larson, J. (2023, May 4). Synthetic Replacements for Human Survey Data? The Perils of Large Language Models. <https://doi.org/10.31235/osf.io/5ecfa>

Kim, J., Byungkyu, L., (2023, Nov 11).
*AI-Augmented Surveys:
 Leveraging Large Language Models and Surveys for Opinion Prediction*
<https://arxiv.org/abs/2305.09620>

DATA: 68,846 individuals' responses to 3,110 questions collected for 33 repeated cross-sectional data between 1972 and 2021 for fine-tuning the LLMs. Retrieved text content of GSS survey questions from GSS data explorer

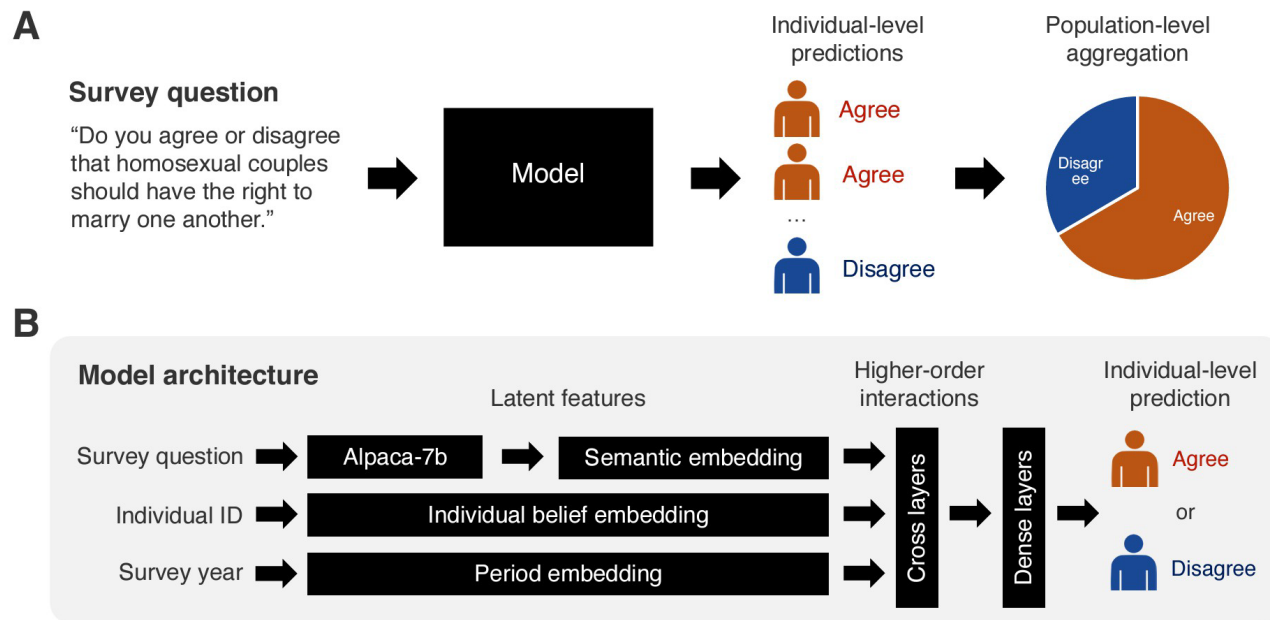


Figure 2: An overview of our methodological framework. In Panel A, we use survey weights when aggregating individual-level prediction into population-level estimates to account for potential sampling bias. In Panel B, individual belief and period embeddings are initially randomly assigned but optimized during the fine-tuning process using dense and cross layers. Semantic embedding, initially estimated by pre-trained LLMs (e.g., Alpaca-7b), is also optimized during the fine-tuning stage.

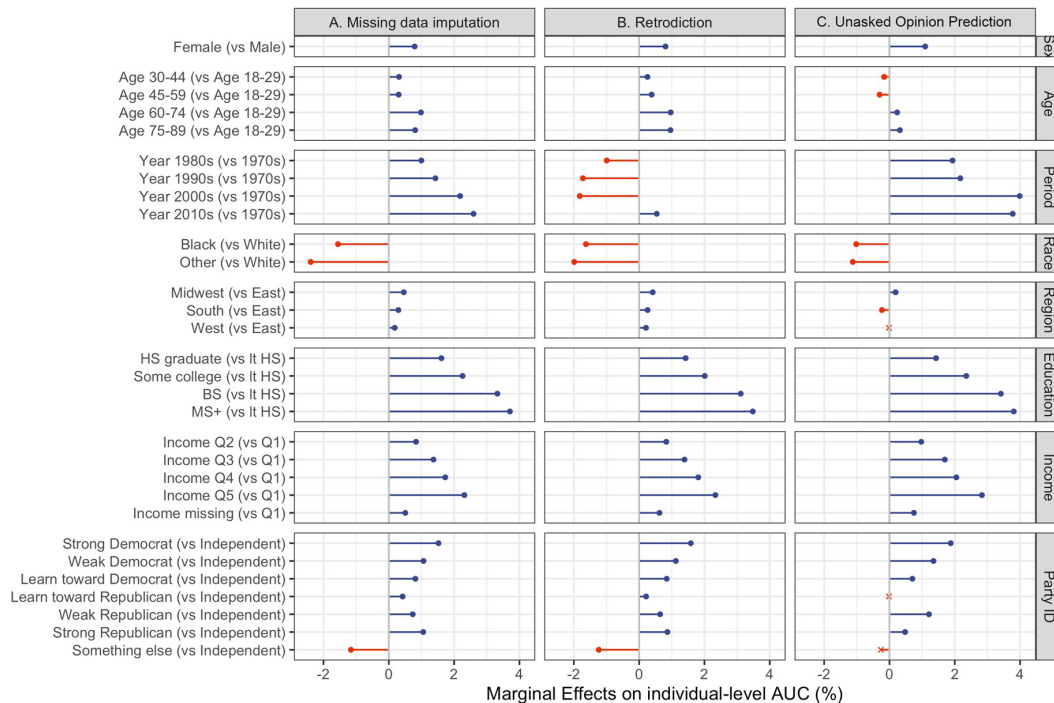
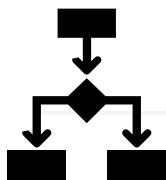


Figure 5: Coefficient plots from OLS regression models predicting individual-level AUC across three different types of missing response prediction. A higher AUC value indicates

For instance, rather than asking the same ten questions to a thousand participants, pollsters can disseminate twenty questions among the same thousand participants, each answering ten questions, and employ the model to **infer individual responses to the remaining ten unasked questions**. On the other hand, given our model's remarkable ability to mimic human responses, even including biases, researchers can use it to **refine their survey questions by systematically examining characteristics of questions that cannot be accurately predicted** (e.g., poor question wording).

Tukey defined data analysis as
“procedures for analyzing data,
techniques for interpreting the results of such procedures,
ways of planning the gathering of data to make its analysis
easier, more precise or more accurate, and all the
machinery and results of (mathematical) statistics which
apply to analyzing data” (Tukey 1962, p. 2)



occupationMeasurement: A Comprehensive Toolbox for Interactive Occupation Coding in Surveys

Jan Simson¹, Olga Kononykhina¹, and Malte Schierholz¹

¹ Department of Statistics, Ludwig-Maximilians-Universität München, Germany Corresponding author

DOI: [10.21105/joss.05505](https://doi.org/10.21105/joss.05505)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Chris Vernon](#)

Reviewers:

- [@welch16](#)
- [@danielruss](#)

Submitted: 30 March 2023

Published: 24 August 2023

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

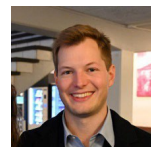
Summary

People earn a living a multitude of ways which is why the occupations they pursue are almost as diverse as people themselves. This makes quantitative analyses of free-text occupational responses from surveys hard to impossible, especially since people may refer to the same occupations with different terms. To address this problem, a variety of different classifications have been developed, such as the International Standard Classification of Occupations 2008 (ISCO) (ILO, 2012) and the German Klassifikation der Berufe 2010 (KldB) (Bundesagentur für Arbeit, 2011), narrowing down the amount of occupation categories into more manageable numbers in the mid hundreds to low thousands and introducing a hierarchical ordering of categories. This leads to a different problem, however: Coding occupations into these standardized categories is usually expensive, time-intensive and plagued by issues of reliability.

Here we present a new instrument that implements a faster, more convenient and interactive occupation coding workflow where respondents are included in the coding process. Based on the respondent's answer, a novel machine learning algorithm generates a list of suggested occupational categories from the Auxiliary Classification of Occupations (Schierholz, 2018), from which one is chosen by the respondent (see Figure 1). Issues of ambiguity within occupational categories are addressed through clarifying follow-up questions. We provide a comprehensive toolbox including anonymized German training data and pre-trained models without raising privacy issues, something not possible yet with other algorithms due to the difficulties of anonymizing free-text data.

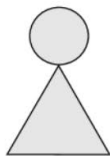
Statement of Need

Assigning occupations to standardized codes is a critical task frequently encountered in research, public administration and beyond: They are used in government censuses (e.g. USA, UK, Germany) and administrative data to better understand economic activity, in epidemiology to estimate exposure to health hazards, and in sociology to obtain a person's socio-economic



ISCO-08: 5141
KldB (2010): 823

The Challenge



You are a software developer who writes python codes from a given set of specifications of the software. Following conditions should be strictly met:

- 1) Every class/method should be fully implemented.
- 2) There should not be any 'pass' statements.
- 3) The packages required should be listed at the end.

Given the following specification, help me write a software in Python that simulates a double pendulum:

list of specifications

Below is the Python code implementation based on the given specifications. Please note that this code is a high-level representation and does not include the actual physics calculations or visualization implementations. Those would require more detailed algorithms and possibly third-party libraries.

Partial codes with placeholder classes



(a) Fails to follow user preference



It is throwing the following error: *Error message*

It appears that there's an incomplete line in the code I provided earlier. The error you're encountering is due to the fact that the code is trying to add a `float` and a `Namespace` object, which is not a valid operation.

Provides new code



The code completed without any error, but no simulation was generated.

I apologize for the oversight. The code provided earlier sets up the simulation and solves the equations of motion but does not include the part that generates the animation.

Provides new code



(b) Unreliable refinement of modular & complex artifact



class Visualization:

```
...
def displayPendulumMotion(self, data_storage):
    # Unpack data
    time_steps = data_storage.time_steps
    angles = data_storage.angles
    ...
```

```
def animate(i):
    # Assuming the length of
    # both arms is 1 for simplicity
    x1 = 1 * np.sin(angles[i][0])
    y1 = -1 * np.cos(angles[i][0])
    x2 = x1 + 1 * np.sin(angles[i][1])
    y2 = y1 - 1 * np.cos(angles[i][1])
    ...
```

Human-edits are just another context!



```
l1, l2 = data_storage.arm_lengths
def animate(i):
    # Assuming the length of
    both arms is 1 for simplicity
    x1 = l1 * np.sin(angles[i][0])
    y1 = -l1 * np.cos(angles[i][0])
    x2 = x1 + l2 * np.sin(angles[i][1])
    y2 = y1 - l2 * np.cos(angles[i][1])
```

(c) Unaware of active human participation

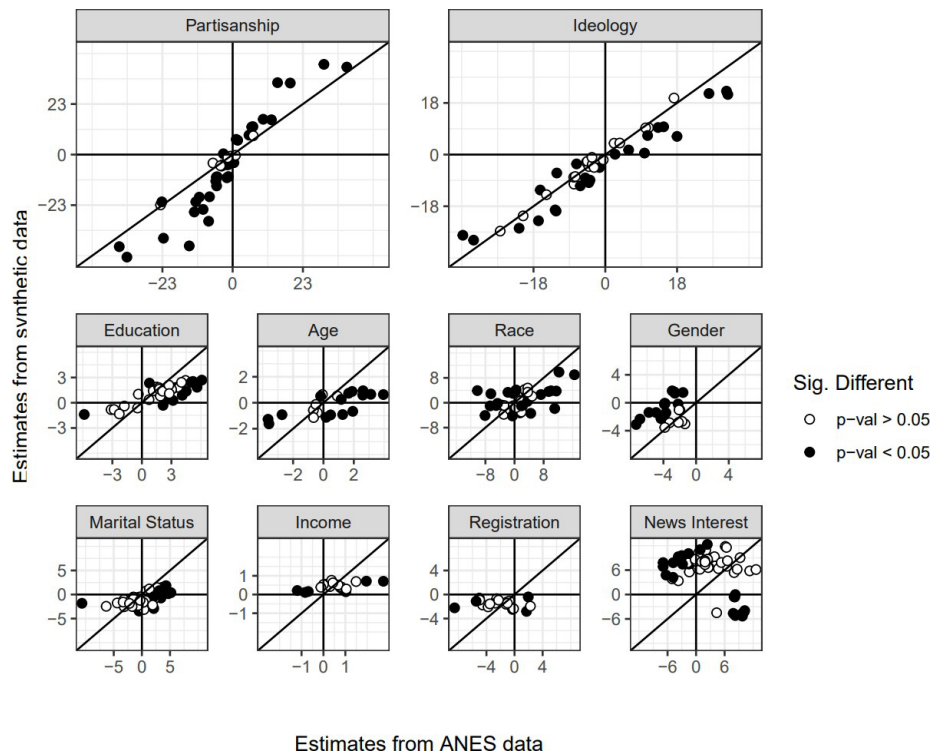


Figure 3: Each point describes the coefficient estimate capturing the partial correlation between a covariate and a feeling thermometer score toward one of the target groups, estimated in either 2016 or 2020. The x-axis position is the coefficient estimated in the ANES data, and the y-axis position is the same coefficient estimated in the synthetic data. Solid points indicate coefficients who are significantly different when estimated in either the ANES or synthetic data, while hollow points are coefficients that are not significantly different. Points in the northeast and southwest quadrants generate the same substantive interpretations, while those in the northwest and southeast quadrants produce differing interpretations. A synthetic dataset that is able to perfectly recover relationships estimated in the ANES data would have all points falling along the 45 degree line.

Bisbee, J., Clinton, J., Dorff, C., Kenkel, B., & Larson, J. (2023, May 4). Synthetic Replacements for Human Survey Data? The Perils of Large Language Models. <https://doi.org/10.31235/osf.io/5ecfa>

English (translation) I am 28 years old and female. I have a college degree, a medium monthly net household income, and am working. I am not religious. Ideologically, I am leaning center-left. I rather weakly identify with the Green party. I live in West Germany. I think the government should facilitate immigration and take measures to reduce income disparities. Did I vote in the 2017 German parliamentary elections and if so, which party did I vote for? I [INSERT]

Notes: We decided not to include “gewählt” (voted) as a suffix in the prompt, using the [MASK] instead of [INSERT] request, as it might bias the output against non-voters by reducing the likelihood of GPT completing the sentence with “nicht” (not) or “ungültig” (invalid) due to German semantics. We leave the further exploration of these effects to prompt engineering researchers.

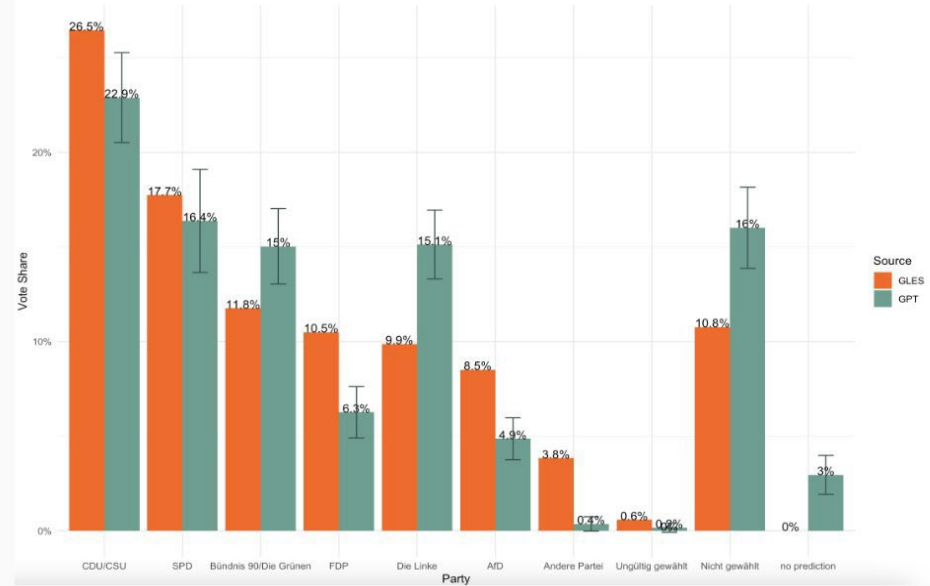
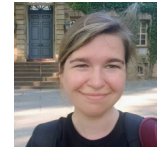


Figure 3: Replicating Argyle et al. for German data (GLES): Current project by Leah von der Heyde, Alexander Wenz and Carolina Haensch

22

Von der Heyde, L., Wenz, A., & Haensch, A.-C. (2024, February 22). Artificial Intelligence, Unbiased Opinions? Assessing GPT’s suitability for estimating public opinion in multi-party systems. <https://doi.org/10.17605/OSF.IO/5BRXD>



The Possible Solution

ExAIC Framework

The ExAIC framework enables coconstruction of solutions
by an expert and a dynamically adaptive agent
through search in a construction space
via natural language communication
with agent integrity by design.

ExAIC Framework

The ExAIC framework enables coconstruction of solutions
by an expert and a dynamically adaptive agent
through search in a construction space
via natural language communication
with agent integrity by design.

ExAIC Framework

The ExAIC framework enables coconstruction of solutions
by an expert and a dynamically adaptive agent
through search in a construction space
via natural language communication
with agent integrity by design.

ExAIC Framework

The ExAIC framework enables coconstruction of solutions
by an expert and a dynamically adaptive agent
through search in a construction space
via natural language communication
with agent integrity by design.

ExAIC Framework

The ExAIC framework enables coconstruction of solutions
by an expert and a dynamically adaptive agent
through search in a construction space
via natural language communication
with agent integrity by design.

ExAIC Framework

The ExAIC framework enables coconstruction of solutions
by an expert and a dynamically adaptive agent
through search in a construction space
via natural language communication
with agent integrity by design.

DS: Real world example

Inspired by the ongoing collaboration of several ExA/CPLs w/ Deutsche Bundesbank



LEGAL ENTITY VALIDATION



DS: Preliminaries



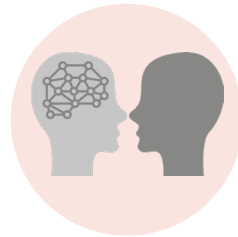
Expert
Data scientist
(= domain expert)



Task
Data analysis and validation



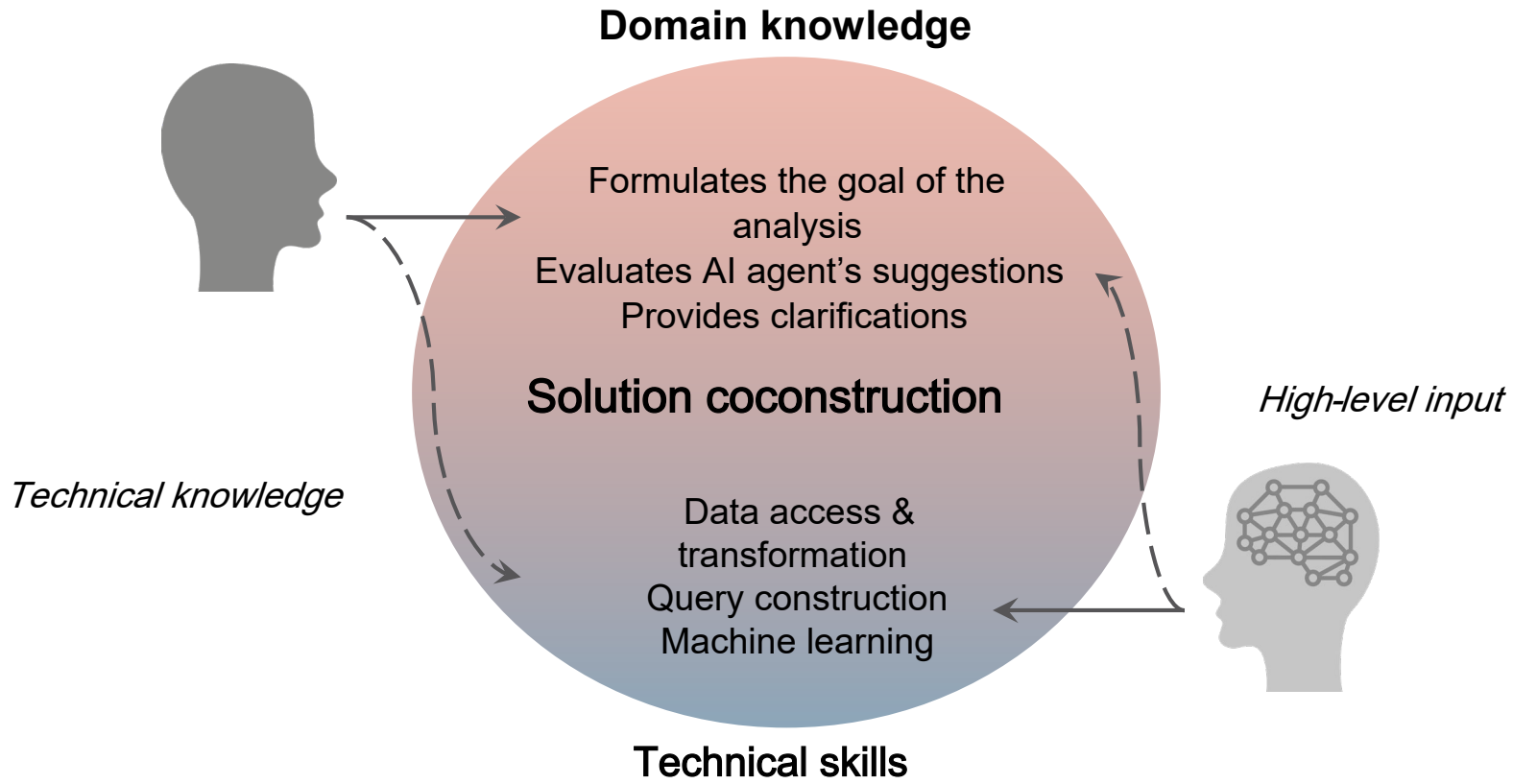
Solution
A data science workflow
(aka *a reasoning chain*) over
several multimodal data sources



Team
Expert and AI *interactively* query,
transform and interpret the data



DS: Expert-AI team





DS: Legal entity validation scenario

Data: private repositories of Deutsche Bundesbank

- **Huge** amounts of data
- Potentially **inconsistent**

Account Holder Name	Company Registration Nr	MS Registry	Installation ID	Installation Name	Activity Type	Permit ID	Permit Expiry	Contact Country	Contact City	Contact PCode	Contact Address L1
Verallia Deutschland AG	HRB 610192	DE	217	Glass production facility Essen	Manufacture of glass	14250-0032	ACTIVE	DE	Bad Wurzach	88410	Oberlandstraße 1-8
Verallia Deutschland AG	HRB 610192	DE	218	Glass production facility Wirges	Manufacture of glass	14250-0033	ACTIVE	DE	Bad Wurzach	88410	Oberlandstraße 1-8
Verallia Deutschland AG	HRB 610192	DE	219	Glass production facility Neuburg	Manufacture of glass	14250-0034	ACTIVE	DE	Bad Wurzach	88410	Oberlandstraße 1-8
Verallia Deutschland AG	HRB 610192	DE	220	Glass production facility Wurzach	Manufacture of glass	14250-0035	ACTIVE	DE	Bad Wurzach	88410	Oberlandstraße 1-8

...



DS: Legal entity validation scenario

Goal: verify that each company exists as a legal entity and its legal data is consistent with the available data sources



Expert can not check it all manually



AI can not solve everything



Expert & AI solution coconstruction !



DS: Legal entity validation scenario

High-level task specification (language)

Task: Find inconsistencies

Area: East Germany

Company types: **Type 1, Type 2, ..., Type 10**

Fine-grained
types

Find inconsistencies
in company type data within
east Germany



Please specify which company
types you are interested in

Fine-grained;
Type 1, ..., Type 10





DS: Legal entity validation scenario

Mid-level workflow specification (structured)

W1:

```
1 For all entries:
2 text = LOOKUP(company type)
3 err = SPELLCHECK(text)
4 REPORT(err)
```

W2:

```
1 For all entries:
2 t1 = LOOKUP(company type)
3 addr = LOOKUP(company address)
4 img = LOOKUP_SATELLITE(addr)
5 t2 = ANALYZE_TYPE(img)
6 err = CHECK(t1, t2)
7 REPORT(err)
```

Check
images

Find inconsistencies
in company type data within
east Germany

Please specify which company
types you are interested in

Fine-grained;
Type 1, ..., Type 10

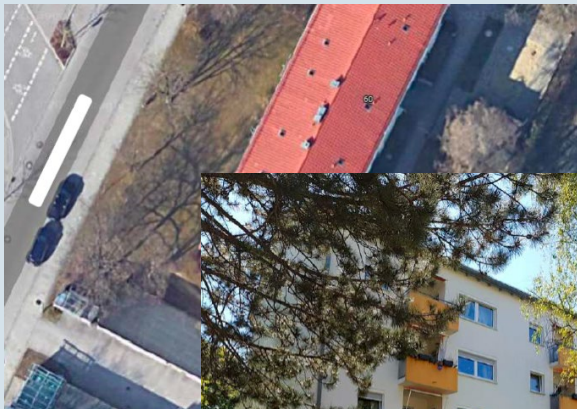
Alright, which workflow do you
prefer, W1 or W2?

W2 is better;



DS: Legal entity validation scenario

**Joinery?
(Manufacturing)**



Find inconsistencies
in company type data within
east Germany



Please specify which company
types you are interested in



Fine-grained;
Type 1, ..., Type 10



Alright, which workflow do you
prefer, W1 or W2?



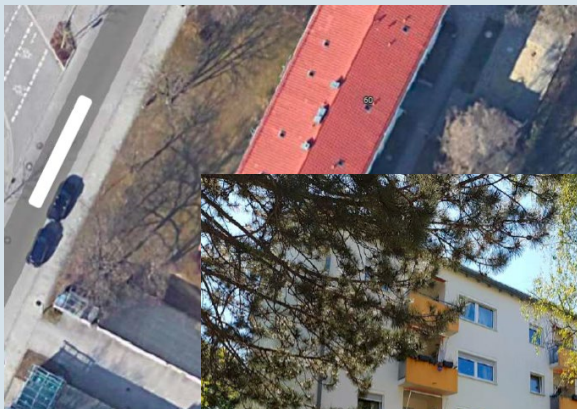
W2 is better;





DS: Legal entity validation scenario

Joinery?
(Manufacturing)



Find inconsistencies
in company type data within
east Germany



Please specify which company
types you are interested in



Fine-grained;
Type 1, ..., Type 10



Alright, which workflow do you
prefer, W1 or W2?



W2 is better;
however, you should also look up
street view images.





DS: Legal entity validation scenario

Mid-level workflow specification (structured)

```
1  For all entries:
2  t1 = LOOKUP(company type)
3  addr = LOOKUP(company address)
4  img1 = LOOKUP_SATELLITE(addr)
5  t2 = ANALYZE_TYPE(img1)
6  img2 = LOOKUP_STREET(addr)
7  t3 = ANALYZE_TYPE(img2)
8  err = CHECK(t1, t2, t3)
9  REPORT(err)
```

Include
street view

Find inconsistencies
in company type data within
east Germany



Please specify which company
types you are interested in



Fine-grained;
Type 1, ..., Type 10



Alright, which workflow do you
prefer, W1 or W2?



W2 is better;
however, you should also look up
street view images.



Ok; executing.





DS: Legal entity validation scenario

Low-level workflow (executable)



Instantiating a solution...

Executing...

157 Errors found

Show me a likely error





DS: Legal entity validation scenario

Production of car
input factors?



Show me a likely error



Here you go; the image shows
a residential building.



With scene view
images, you should also look
left and right.





DS: Legal entity validation scenario

Production of car
input factors?



Show me a likely error



Here you go; the image shows
a residential building.



With scene view
images, you should also look
left and right.



Got it; this was a false positive.





DS: Legal entity validation scenario

Mid-level workflow specification (structured)

```
1   For all entries:
2   t1 = LOOKUP(company type)
3   addr = LOOKUP(company address)
4   img1 = LOOKUP_SATELLITE(addr)
5   t2 = ANALYZE_TYPE(img1)
6   [images] = LOOKUP_STREET(addr)
7   types = ANALYZE_TYPE(images)
8   err = CHECK(t1, t2, types)
9   REPORT(err)
```

Multiple
views

Show me a likely error



Here you go; the image shows
a residential building.



With scene view
images, you should also look
left and right.



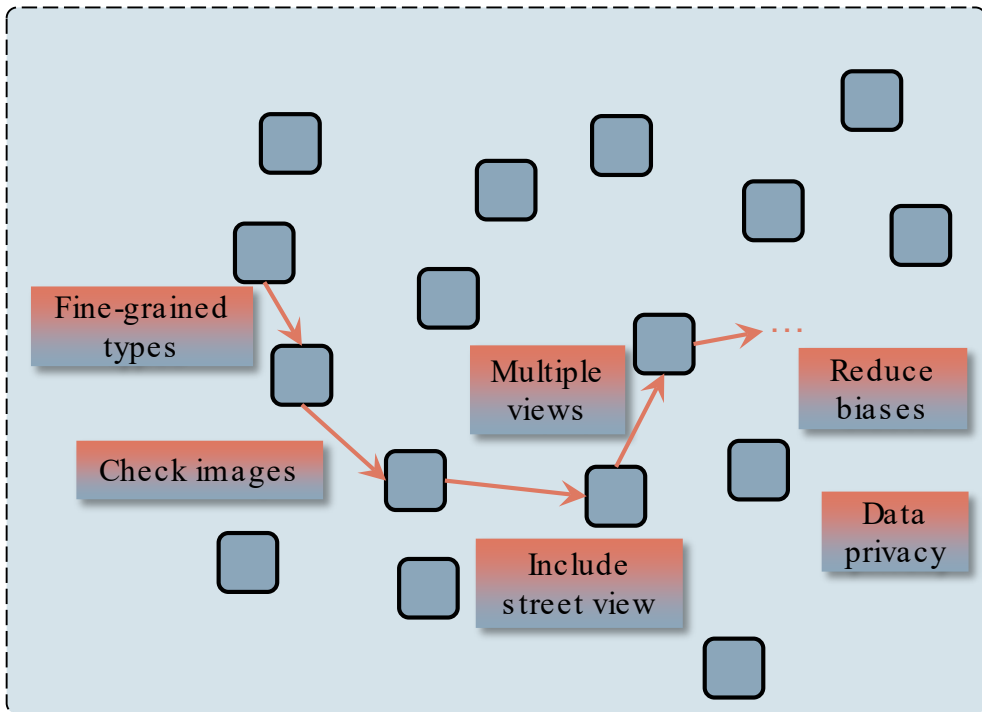
Got it; this was a false positive.





DS: Legal entity validation scenario

Construction space (mid-level workflow)



Over time:

- The agent learns the expert's preferences
- The quality/quantity of communication evolves





AI agents **need to know how to ask** clarification questions so that the expert provides the right answers



This requires **expert-awareness** and generation of the **right clarification questions**



We develop a method that **incrementally learns from sparse** and **multimodal expert inputs**



I have a dataset with information about companies and I need to check whether the industry classifications of the companies are correct. Can you help me with this?

Of course! Do you want to **use the companies with the highest number of employees** for the selective checks?

1



I have a dataset with information about companies and I need to check. Based on a previous check with a different dataset, **I don't expect companies in the construction sector to be misclassified. So, I would leave those out.**

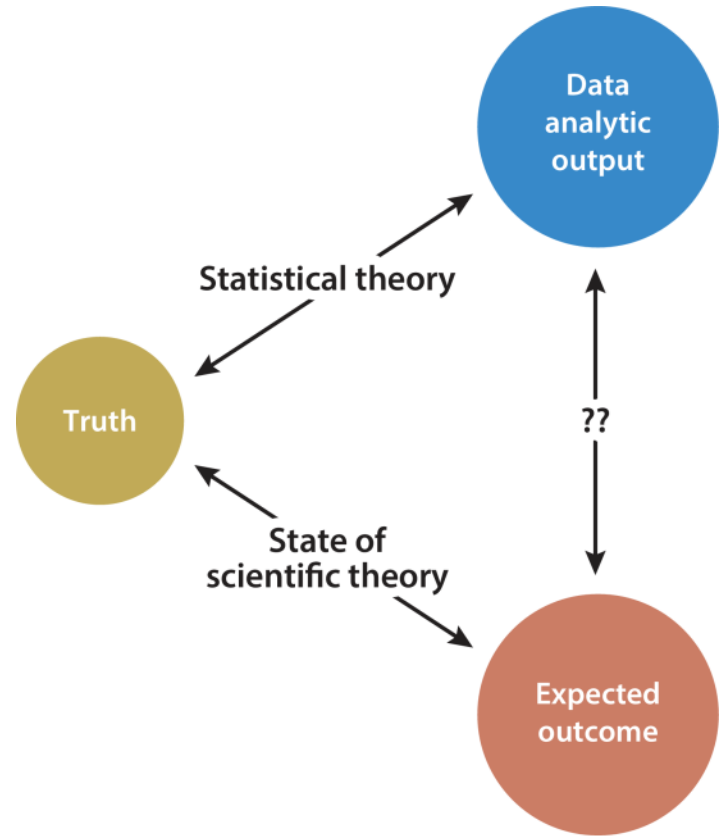
2

Are you sure you want to leave out companies in the construction sector? **Your previous finding might anchor your expectations and does not necessarily apply to the current dataset.**

3

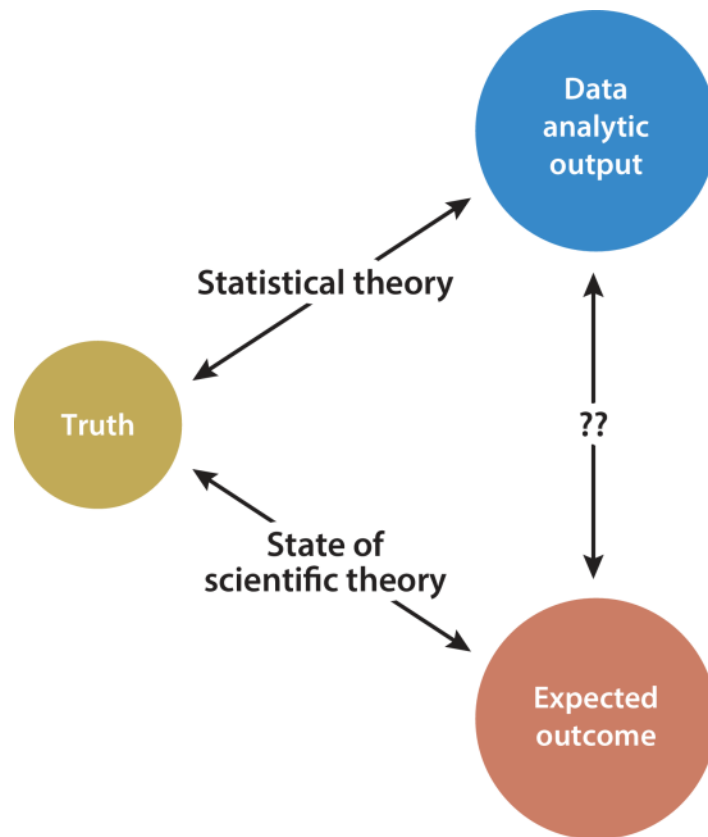


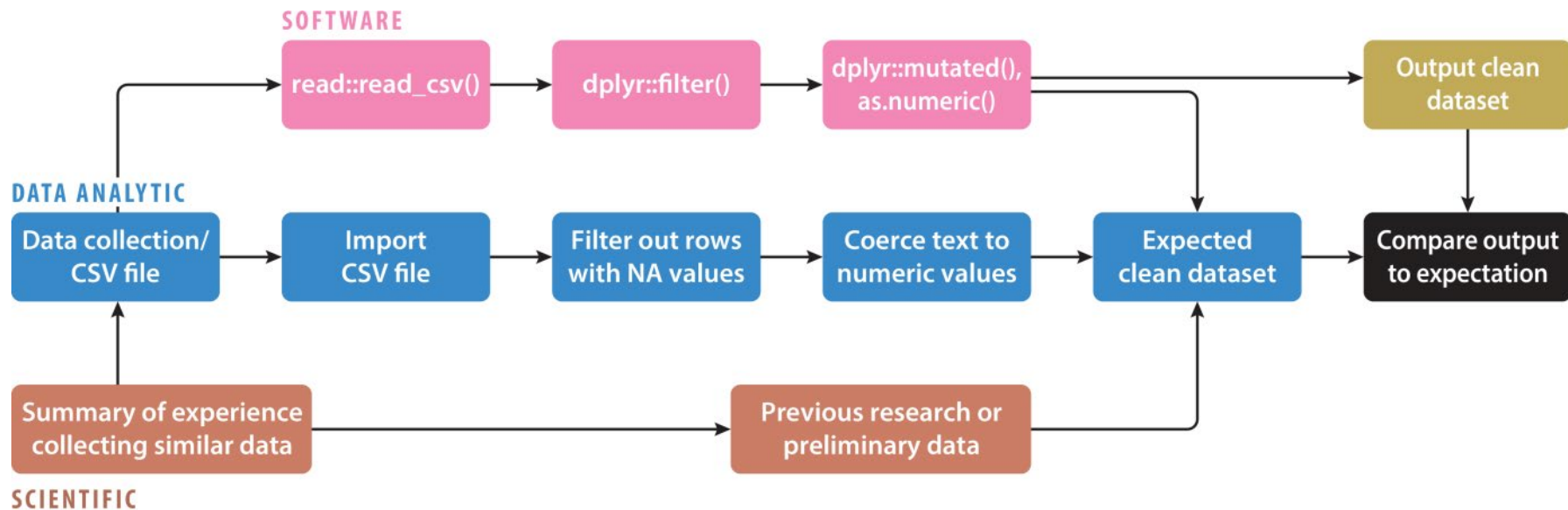
The (missing) Data Science Theory

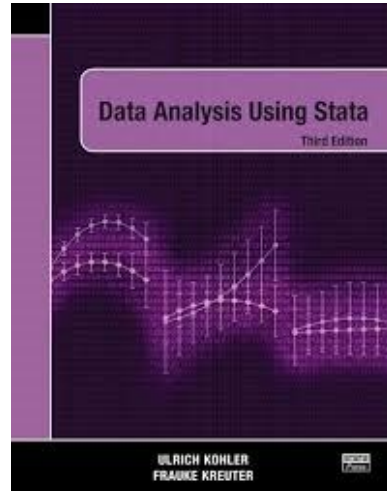
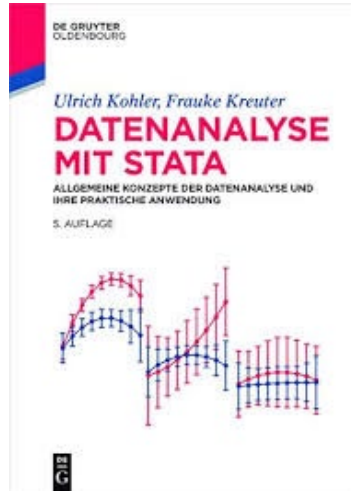


Research Questions

- 1 What data science steps should the agent **be aware of**?
- 2 How can we **enable an agent to ask clarification questions** that address the cognitive biases a human expert may have?
- 3 How can we **evaluate the suitability of the agent-expert dialog efficiently** with sparse data?







Data Analysis With AI Using Stata