

Data analysis tools for statistical non-experts

Eunice Jun

A dissertation
submitted in partial fulfillment of the
requirements for the degree of

Doctor of Philosophy

University of Washington
2023

Reading Committee:
Jeffrey Heer, Co-Chair
René Just, Co-Chair
Tyler H. McCormick

Program Authorized to Offer Degree:
Paul G. Allen School of Computer Science & Engineering

© Copyright 2023

Eunice Jun

University of Washington

Abstract

Data analysis tools for statistical non-experts

Eunice Jun

Co-chairs of the Supervisory Committee:

Jerre D. Noe Endowed Professor Jeffrey Heer
Paul G. Allen School of Computer Science & Engineering

Associate Professor René Just
Paul G. Allen School of Computer Science & Engineering

Data analysis is critical to science, public policy, and business. Despite their importance, statistical analyses are difficult to author, especially for researchers with expertise outside of statistics. Existing statistical tools, prioritizing mathematical expressivity and computational control, are low-level while researchers' motivating questions and hypotheses are high-level. Researchers need to translate their questions and hypotheses into low-level statistical code in an error-prone process that involves grappling with their domain knowledge, statistics, and programming.

This dissertation introduces two tools that embody a new way of authoring analyses: Tea and Tisane. Researchers directly express their domain knowledge through higher level abstractions, and the tools will validate the data, select a statistical analysis, and implement it, all while educating analysts about why a statistical approach is valid. Tea helps analysts author statistical tests. Tea's key insight is that statistical test selection can be cast as a constraint satisfaction problem. Tisane enables analysts to author generalized linear models with or without mixed effects, which are difficult for even statistical experts to author. Using Tisane, analysts can express their conceptual models using a high-level domain specific language. Tisane translates these conceptual models into causal DAGs and engages analysts in a disambiguation process to arrive at an output statistical model. Real-world researchers have already used these tools to conduct analyses in published research that push their own disciplines forward. I will also introduce "hypothesis formalization," a series of cognitive and operational steps analysts take to translate their research questions into statistical implementations. Hypothesis formalization retrospectively explains why Tea improves statistical testing and directly inspired the design of Tisane.

Tea and Tisane serve as platforms for further research into computational support for statistical analysis. This work also exemplifies how combining human-computer interaction with other areas

in and outside of computer science leads to software tools that impact real-world users.

Acknowledgements

I am lucky to have many, many people to thank.

I thank my advisors Jeffrey Heer and René Just. Jeff, thank you for your endless creativity, sense of humor, and focused attention in all our conversations. I had a lot of fun and freedom during the second half of my PhD thanks to you. Thanks for helping me navigate personal and professional challenges with integrity. René, thank you for your consistent insistence on the details.

I have the best committee, hands down. Emery Berger, I am deeply grateful that you agreed to meet with me, then a total stranger, at the start of your sabbatical at MSR in 2018. Who knew that coffee would lead to many, many more espressos while we worked on Tea, Scone, and many other projects. You have been my most ardent champion and creative conversational partner. Thank you for including me in your latest crazy idea and laughing with me through the highs and lows of graduate school. Leilani Battle, thank you for your continual enthusiasm and support for my research and mentorship since you were a postdoc at UW! It's great to have you back. Tyler McCormick, thank you for welcoming me into your research group and the Center for Statistics and the Social Sciences (CSSS). My thinking and scholarship have expanded due to our conversations and collaboration.

I feel lucky to have been part of the warm UW CSE community throughout my PhD. In particular, I am indebted to the members of the Interactive Data Lab (IDL) with whom I overlapped during my PhD for their camraderie: Arvind Satyanarayan, Ham Wongsuphasawat, Dominik Moritz, Yeaseul Kim, Younghoon Kim, Jane Hoffswell, Sherry Wu, Alex Kale, Yang Liu, Matt Conlen, Zening Qu, Mick Kittivorawong, Ameya Patil, Junran Yang, Madeleine Grunde-McLaughlin, Josh Horowitz, and Luke Snyder. Michael Correll and Maureen Stone, while not students in IDL, were consistent members of the lab. Thank you, all, for keeping me humble, honest, and hopeful about the positive social impact of great research and open-source software.

I have also been fortunate to be part of the idiosyncratic PLSE lab. To my PLSE labmates over the years, of which there are too many to name individually, I love how ferociously curious and whimsical you are about everything in research and life. I would not have pursued this research without PLSE welcoming me early on. In this regard, I want to give special thanks to my friends Sarah Chasins, James Bornholt, Sami Davies, Chandra Nandi, Amanda Swearngin, Chenglong Wang, Mangpo Phothilimthana, Jasper Tran O'Leary, Pavel Panchevka, and Max Willsey. Zachary Tatlock, thank you for investing so much of yourself into making PLSE the best place to pursue daring research, have fun, and make friends.

To Elise Dorough, thank you for everything you do to keep CSE running while still making time to help my peers and me through any situation (and many emotions). To Joe Eckert, thank you for answering my

numerous Slack messages and cries for help even when I caused the logistical problems myself. To Elle Brown, thank you for your kindness. To Lisa Merlin, thank you for helping me with countless reimbursements and doing everything in your power to make research smooth.

The larger UW ecosystem has also shaped my scholarship. To the members of CSSS, thank you for welcoming me and giving me many opportunities to learn from and ask questions of scholars across the social and statistical sciences. The interdisciplinary seminars and courses, such as Thomas Richardson's offering of CSSS/STAT 566 on causal reasoning, were instrumental in my graduate education. Thanks also to the larger HCI community of DUB. There are a few individual mentors to whom my gratitude is overdue: Daniela Rosner (HCDE), Tim Althoff (CSE), and Amy Ko (iSchool). To Daniela, especially, thank you for sharing your brilliance, wisdom, and curiosity about "design" and life over memorable chats in your office and on walks throughout the city.

Throughout my PhD, I was also fortunate to expand my research horizons through two summer internships at Microsoft Research with Mary Czerwinski, Daniel McDuff, and Ben Zorn. Thank you for entrusting me to run with new ideas, modeling constructive mentorship styles, and becoming friends over the years.

To my collaborators, thank you, for your inspiration, energy, and support. Specifically, I thank Joseph Dieleman and Goli Tsakalos, for inviting me as a volunteer with the Institute for Health Metrics and Evaluation (IHME) the last two years of my PhD. To Sawyer Crosby and Emily Johnson, thank you for your time and energy in showing me the ropes at IHME, sharing your analysis experiences, and trying out and giving helpful feedback on early prototypes of my systems. To Maureen Daum, thank you for your candor, reliability, and sense of humor as a collaborator and friend. To Nicole de Moura, Grace Oh, Melissa Birchfield, Pranav Rajan, Shreyash Nigam, Josh Pollock, Annie Denton, Blue A. Jo, Reiden Chea, Corinne Herzog, Vincent Pun, Irene Luo, Ken Gu, and Audrey Seo, thanks for entrusting me to mentor you during your explorations of research and computer science. I am so immensely proud of how you have pursued your interests, and I look forward to continue seeing you grow!

I am immensely grateful to my community in and around Seattle. To my Council of Wise Women from Advent Anglican, Seattle Mosaic Arts, and St. Mark's Gardening Ministry, thank you for enriching my soul. You have shared so much wisdom in the countless hours I have sat in your presence, often in silence.

I give my heartiest, warmest thank you to "The Emotional B*tches (Thoughts-Emotions-Behaviors)" and the healthcare professionals at The Emily Program. You know who you are. My heart is beyond a 5 on the fullness scale. It's overflowing with marbles. Thank you for creating a true fellowship and, quite literally, keeping me alive. Patrice Staiger and Carrie DeMartini, thank you for guiding me to and through The Emily Program. Treg Isaacson, Hui Sun, and the others at Hall Health, thank you for lighting the way.

I thank my friends and family. To Bridget Claborn, thank you for your incisive questions, belief in the poetry and beauty around us (i.e., post-post-modernism, dahlias, anemones, Mary Oliver), and commitment to weave together research and practice for social change. To Aoife Blacklaws, thank you for encouraging me in all my endeavors and sharing your awe of every little creature. To Michael Zuch, thank you for sharing in the tender questions of faith, theology, and the meaning of life. To Rachael Grenfell-Dexter, thank you for sharing many adventures (all starting in Guatemala!) and giving me permission to have fun, grow and evolve. To Lucy Rahner, thank you for sharing your artistry and earnest search for the truth and goodness in everything. To Catherine Chung, thank you for keeping it real, talking with me endlessly for days, and

listening so whole-heartedly. To the friends named and unnamed, I am honored to share this life with you.

Finally, to Jared Roesch, thank you for your love, commitment to community, and ambition to focus on and fight for what truly matters. To my parents, thank you for your sacrifices and support. I love that you encourage me to write better and more papers even when, in your estimation, no one reads them.

DEDICATION

To Bill Thompson:

Thank you for sharing your curiosity, sharp intellect, patience, and kindness.

Thanks for modeling how to pursue important scientific endeavors and bring people along.

You're with me always as I design studies and interpret data, crack open a can of beer at the top of a mountain, and orienteer through the world and life.

Cheers to you.

Contents

| | | |
|-----------------------------|--|-----------|
| 1 | Introduction | 19 |
| 1.1 | Thesis Approach and Statement | 19 |
| Thesis statement | 19 | |
| Challenge 1 | 20 | |
| Challenge 2 | 20 | |
| Challenge 3 | 20 | |
| 1.2 | Summary of Contributions | 21 |
| 1.3 | Thesis outline | 22 |
| 2 | Related work | 25 |
| 2.1 | Statistical data analysis as sensemaking | 25 |
| 2.2 | Empirical accounts of data analysis practice | 26 |
| 2.3 | Tools for data analysis | 27 |
| 2.3.1 | Tools for conceptual modeling | 28 |
| 2.3.2 | Tools for study design | 28 |
| 2.3.3 | Tools for statistical specification | 29 |
| 2.4 | Validity in statistical data analysis | 29 |
| 3 | Tea: A Domain-Specific Language and Runtime System for Hypothesis Testing | 31 |
| 3.1 | Background and Related work | 32 |
| Statistical scope | 33 | |
| 3.2 | Usage Scenario | 34 |
| 3.3 | Design Considerations | 35 |
| 3.4 | System overview | 37 |
| 3.4.1 | Tea’s Domain-Specific Language | 37 |
| 3.4.2 | Tea’s Constraint-based Runtime System | 40 |
| 3.5 | Evaluation | 42 |

| | | |
|---|--|-----------|
| 3.5.1 | How does Tea compare to textbook tutorials? | 42 |
| 3.5.2 | Does Tea avoid common mistakes made by non-expert users? | 44 |
| 3.6 | Discussion, Limitations, and Future Work | 44 |
| 3.7 | Summary of Contributions | 45 |
| 4 | Hypothesis Formalization: A conceptual framework describing how analysts translate research questions into statistical analyses | 49 |
| 4.1 | Background and Related Work | 52 |
| 4.1.1 | Statistical Thinking | 52 |
| 4.1.2 | Statistical data analysis as part of scientific discovery | 53 |
| 4.2 | Formative content analysis | 54 |
| 4.2.1 | Discussion | 58 |
| 4.2.2 | Takeaways: Expected Steps in Hypothesis Formalization | 59 |
| 4.3 | Exploratory Lab Study | 60 |
| 4.3.1 | Methods | 60 |
| 4.3.2 | Findings and Discussion | 62 |
| 4.3.3 | Takeaways from the Lab Study | 69 |
| 4.4 | Analysis of Software Tools | 70 |
| 4.4.1 | Method | 70 |
| 4.4.2 | Findings and Discussion | 71 |
| 4.4.3 | Takeaways from the Analysis of Tools | 76 |
| 4.5 | Design Implications for Statistical Analysis Software | 76 |
| Connect Model Implementations with Mathematical Equations | 76 | |
| Express Conceptual Hypotheses to Bootstrap Statistical Model Implementation | 77 | |
| Co-author Conceptual and Statistical Models | 78 | |
| 4.6 | Discussion | 79 |
| 4.7 | Future Work | 80 |
| 4.8 | Summary of Contributions | 81 |
| 5 | Tisane: Authoring statistical models via formal reasoning from conceptual and data relationships | 83 |
| 5.1 | Background and Related work | 84 |
| 5.1.1 | Statistical scope | 85 |
| 5.2 | Early Design Process | 86 |
| 5.3 | First Release | 87 |
| 5.3.1 | Study design specification language and graph representation | 88 |

| | | |
|----------|--|------------|
| 5.3.2 | Statistical model inference: Interactively querying the graph IR | 93 |
| 5.4 | Initial evaluation: Case studies with researchers | 99 |
| 5.4.1 | Case Study 1: Planning a new study | 99 |
| 5.4.2 | Case Study 2: Analyzing data for a paper submission | 100 |
| 5.4.3 | Case Study 3: Developing models to inform future models | 101 |
| 5.5 | Limitations and Motivation for Re-design | 103 |
| 5.6 | Elicitation lab study | 103 |
| 5.6.1 | Method | 104 |
| 5.6.2 | Key Observations | 104 |
| 5.6.3 | DSL Re-design Goals | 107 |
| 5.7 | Second Release: rTisane | 107 |
| 5.7.1 | rTisane’s DSL | 109 |
| 5.7.2 | Two-step interactive compilation | 110 |
| 5.8 | Summative Evaluation: Controlled lab study | 113 |
| 5.8.1 | Study design | 113 |
| 5.8.2 | Analysis Approach | 114 |
| 5.8.3 | Findings | 115 |
| 5.8.4 | Discussion | 121 |
| 5.8.5 | Limitations and Future Work on rTisane | 122 |
| 5.9 | Discussion | 123 |
| 5.10 | Summary of Contributions | 123 |
| 6 | Conclusion | 125 |
| 6.1 | Discussion | 125 |
| 6.1.1 | Challenge 1: Designing the <i>right</i> level of abstraction | 126 |
| 6.1.2 | Challenge 2: Representing and reasoning about analysis decisions | 126 |
| 6.1.3 | Challenge 3: Interaction as reflection | 127 |
| 6.2 | Recent developments | 127 |
| 6.2.1 | Construct validity: Within reach with the usage of LLMs | 127 |
| 6.2.2 | What about in the face of LLMs? | 127 |
| 6.3 | Limitations and Future work | 127 |
| 6.3.1 | Support interpretation of statistical results | 127 |
| 6.3.2 | Connect statistical modeling and testing | 128 |
| 6.3.3 | Develop a grammar of study design | 128 |
| 6.3.4 | Support more phases of the data lifecycle | 128 |

| | | |
|--|---|------------|
| 6.3.5 | Promote analytical best practices in science | 128 |
| 6.3.6 | Improving data science education | 129 |
| 6.4 | Impact | 130 |
| 6.5 | Closing Remarks | 130 |
| Bibliography | | 132 |
| A Appendix A: Content analysis resources | | 151 |
| A.1 | Dataset overview | 151 |
| A.2 | Procedure | 151 |
| A.3 | Codebook | 152 |
| A.4 | Additional findings: Contribution types | 152 |
| A.5 | Summary of papers analyzed | 154 |
| B Appendix B: Tisane’s First Release: Additional Examples | | 161 |
| B.1 | Additional examples of graphs that may be constructed | 161 |
| B.2 | Cautioning analysts about adding certain kinds of variables | 162 |
| C Appendix C: Exploratory study materials | | 163 |
| C.1 | rTisane study materials | 163 |

List of Figures

| | | |
|-----|--|----|
| 3.1 | Sample Tea program. The specification outlines an experiment to analyze the relationship between geographic location ('So') and probability of imprisonment ('Prob') in a common USCrime data set [VR13; Kab11]. See Section 3.2 for an explanation of the code. Tea programs specify 1) data, 2) variables, 3) study design, 4) assumptions, and 5) hypotheses. | 34 |
| 3.2 | Tea program and its mode-dependent executions. a) Tea program that aims to determine if two contributor variables, 'Illiteracy' and 'HS Grad' that may predict a third outcome variable 'Life Exp', are correlated. The user asserts that 'Illiteracy' is normally distributed. b) By default, Tea executes programs in the <i>strict</i> mode. c) Warning that Tea disagrees with the user and will override the user's assertion that 'Illiteracy' is normally distributed in the <i>strict</i> mode. d) Results without the parametric test since Tea overrides user's assertion. e) A single line change can modify Tea to execute a program in <i>relaxed</i> mode. f) Warning that Tea cannot verify normality for 'Illiteracy' but will defer to user's assertion. g) Results with the parametric test since Tea proceeds as if 'Illiteracy' was normally distributed. | 47 |
| 4.1 | Definition and overview of the hypothesis formalization steps and process. | 51 |
| 4.2 | Relationship between hypothesis formalization and prior work. | 54 |
| 4.3 | Formative content analysis: example reorderable matrix for [NHNO19]. | 55 |
| 4.4 | Sample statistical specification (D8). | 64 |
| 5.1 | The graph representation of the variables and relationships from the usage scenario. causes edges are labeled with "causes". associates_with edges are labeled with "assoc." Dashed edges indicate nests_within relationships, and dotted edges indicate has relationships. | 89 |

| | |
|---|-----|
| 5.2 Example Tisane GUI for disambiguation. Tisane asks analysts disambiguating questions about variables that are conceptually relevant and that analysts may have overlooked in their query. (A) The left hand panel gives an overview of the model the analyst is constructing. (B) Based on the variable relationships analysts specify, Tisane infers candidate main effects that may be potential confounders. Tisane asks analysts if they would like to include these variables, explaining in a tooltip (C) why the variable may be important to include. (D) Tisane only suggests interaction effects if analysts specify moderating relationships in their specification. This way, Tisane ensures that model structures are conceptually justifiable. (E) From the data measurement relationships analysts provide, Tisane automatically infers and includes random effects to increase generalizability and external validity of statistical findings. (F) Tisane assists analysts in choosing an initial family and link function by asking them a series of questions about their dependent (e.g., Is the variable continuous or about count data?). To help analysts answer these questions and verify their assumptions about the data, Tisane shows a histogram of the dependent variable. | 97 |
| 5.3 rTisane’s conceptual model disambiguation interface. A. Options for resolving ambiguities in the conceptual model due to <i>relates</i> relationships. B. Check and follow-up questions for breaking any cycles that hinder statistical model derivation. C. Graph visualizing expressed conceptual model. | 111 |
| 5.4 rTisane’s statistical model disambiguation interface. | 112 |
| 5.5 Example conceptual models from participants in the summative evaluation without using rTisane. | 115 |
| B.1 More complex examples of <i>moderates</i> written in Tisane’s study design specification language, and their representation in Tisane’s graph IR. Variables are named with u for units, m for measures, and v for data variables that can be either units or measures. Black edges have been added due to the <i>moderates</i> relationship. Gray edges already existed in the graph. In (a), only m1 is a measure, whose unit is u2, so u1*m1 inherits an attribution edge only from u2. In (b), m1 and m2 are measures, with units u1 and u2 respectively, so m1*m2 inherits attribution edges from both u1 and u2. In (c), measures m1 and m2 share a unit, u, and m1*m2 inherits only one attribution edge from u. | 161 |

- B.2 A graph demonstrating an edge case for candidate main effect identification, where the graph contains only associative edges. Candidate main effects are labeled “CME”, independent variables “IV”, and dependent variables “DV”. Variables that are none of the above are left unlabeled. When a graph contains only associative edges, candidate main effects are identified as those that are either associated with the DV or are associated with both the IV and the DV. (Note that the graph could contain additional edges/nodes other than the ones pictured, but the additional edges would not violate any of the initial checks that Tisane makes on the graph IR.) 162
- B.3 An example of the warning text given for potential confounding associations. When analysts hover over the “Warning” badge, a tooltip pops up that explains that they should be careful about adding this variable. Associative relationships may in actuality be causal relationships, and if in fact `pounds_lost caused age`, then adding `age` would invalidate the model. 162

List of Tables

| | | |
|-----|---|-----|
| 3.1 | Comparison of Tea to other tools. Despite the published best practices for statistical analyses, most tools do not help users select appropriate tests. Tea not only addresses the best practices but also supports reproducing analyses. | 36 |
| 3.2 | Results of applying Tea to 12 textbook tutorials. | 43 |
| 4.1 | Overview of the software tools included in our analysis. | 72 |
| 5.1 | Overview of study design tools that informed Tisane’s study design specification language. The first five tools provide higher-level abstractions. They are designed to help researchers reason about their study designs more holistically. The latter eight tools are lower-level and are more focused on stimuli, trials, and progressions between trials. *JsPsych is the base package to which JsPsychR, xprmtnr, and Jaysire provide wrappers and extensions. | 88 |
| 5.2 | The available family and link functions in Tisane. Tisane generates code to fit models using statsmodels and pymer4. The package statsmodels supports GLMs without mixed-effects and a wider variety of family and link function combinations. The package pymer4 supports GLMs with mixed effects and has much more limited support for family and link functions. As statsmodels and pymer add more support, Tisane can be extended. | 96 |
| 5.3 | Participants in summative evaluation. | 114 |
| A.1 | The codebook for analyzing the content of research publications. | 153 |
| A.2 | Summary of CHI papers in our dataset. | 155 |
| A.3 | Summary of JFE papers in our dataset. | 156 |
| A.4 | Summary of Nature papers in our dataset. | 157 |
| A.5 | Summary of PNAS papers in our dataset. | 158 |
| A.6 | Summary of PS papers in our dataset. | 159 |

Chapter 1

Introduction

Statistical analysis plays a critical role in how people make decisions. Policy makers rely on models to track disease, inform health recommendations, and allocate resources. Scientists develop, evaluate, and compare theories based on data. Journalists report on new findings in science, which individuals use to inform decisions that impact their nutrition, finances, and other aspects of their lives. Faulty statistical models can lead to spurious estimations, findings that do not generalize or reproduce, and a misinformed public.

In the context of scientific research, accurate statistical analyses are essential to scientific reproducibility. In a 2016 Nature survey, 1,500 identified “selective reporting”, “pressure to publish,” and “low statistical power or poor analysis” as the top three contributors to the reproducibility crises in their disciplines [Bak16]. The scientists also articulated that a “better understanding of statistics” would be the best approach to improve reproducibility [Bak16].

Despite the prevalence and central importance of statistical analyses, they remain challenging to author accurately. Key to analysis authoring is grappling with and translating implicit domain knowledge into statistical models executable in code [WP99; CEG⁺16]. This dissertation hypothesizes that better understanding and supporting this translation process can enable statistical non-experts author analyses accurately.

1.1 Thesis Approach and Statement

This dissertation (i) probes into how and why translating domain knowledge into executable statistical models in code is difficult and (ii) develops new computational tools that help statistical non-experts author valid analyses by integrating disciplinary, statistical, and programming details. Moving between building systems and empirically studying analysts, this dissertation demonstrates the following:

Thesis statement A combination of domain-specific languages (DSLs) and automated reasoning can help statistical non-experts more readily author valid analyses. Analysis DSLs express conceptual knowledge, data collection procedures, and analysis intents. Automated reasoning methods then compile the conceptual DSL specifications into statistical analysis code.

The following three challenges fall out of this thesis statement.

Challenge 1: Make implicit domain knowledge explicit

Designing abstractions focused on conceptual knowledge requires identifying what domain knowledge analysts want and can express and then balancing these constraints with what automated reasoning approaches may require. What is easy to express and what is easy to assume for the sake of automation may be at odds, especially when analysts provide ambiguous specifications that could be compiled into multiple statistical analyses. The challenge, therefore, is to design language constructs that are usable for analysts and useful for automated reasoning, leveraging interactive program specification as necessary.

Challenge 2: Represent and reason about key statistical analysis decisions

A central idea in this thesis is that software systems should take on the responsibility of translating conceptual knowledge into statistical analyses. This is akin to representing the conceptual knowledge analysts express and compiling it to statistical analyses that respect statistical best practices and rules. To achieve this, picking representations where the reasoning approach “falls out” is a technical challenge.

Challenge 3: Increase analysts’ statistical understanding

While automating statistical analysis can be helpful, analysts relying on data to make high-impact decisions (e.g., policy, scientific discovery) often need to understand why an analysis approach is appropriate and what the implications of the results are to their domain. Furthermore, software can inform how users approach future analyses. Therefore, educating analysts about the applicability and impact of statistical decisions and guiding their interpretation of results are important.

1.2 Summary of Contributions

This dissertation contributes principles and systems for designing statistical analysis tools for statistical non-experts. The contributions can be summarized as follows:

1. A conceptual framework characterizing statistical analysis authoring.
 - (a) Our theory of *hypothesis formalization* describes the cognitive and operational steps involved in translating a high-level conceptual research question and hypothesis into a statistical analysis implemented in code. *Hypothesis formalization* explains why existing statistical tools fail to support statistical non-experts and informs the design of systems developed in this dissertation.
 - (b) We provide the first account scrutinizing the *hypothesis formalization* process in situ. Whereas previous studies of data analysis have relied on self-reports about analysis processes, we conduct an in-depth lab study where we observe analysts prepare and even start to implement statistical models.
 - (c) We qualitatively assess 20 statistical analysis libraries and standalone systems and illustrate how their designs represent the current ecosystem of statistical tools and presently influence data analysis practice. Furthermore, combining our theory of *hypothesis formalization* and this assessment of tools, we develop three design implications for how data analysis software could better serve statistical non-experts.
2. The design, implementation, and evaluation of new DSLs. These DSLs explore ways to design abstractions that prioritize making implicit domain knowledge explicit.
 - (a) The Tea DSL provides a high-level API so that analysts can make explicit their assumptions about the data and their hypotheses to assess using Null Hypothesis Significance Tests.
 - (b) The Tisane DSL captures analysts’ “fuzzy” assumptions about how variables relate in their discipline. The variables and relationships comprise a *conceptual model*.
 - (c) A formative study showed how statistical non-experts implicitly think about causality, how they would like to express their implicit assumptions, and what they expect language constructs describing conceptual models to mean. These findings informed the re-design of Tisane’s DSL, released as rTisane.
 - (d) A benchmark comparison, series of case studies, and a controlled lab study demonstrate the benefit of these DSLs in helping analysts become more aware of their implicit assumptions.

3. Formal representations and automated reasoning approaches for statistical analysis authoring.

To support statistical testing and modeling, we develop representations that allow automated reasoning to compile conceptual models into statistical models.

- (a) In Tea, we implement a constraint-based model and knowledge base for Null Hypothesis Significance Tests.
- (b) In Tisane, we develop an intermediate graph representation to summarize key conceptual assumptions and data collection details. Importantly, a subgraph of the representation is a causal diagram useful for deriving statistical models formally.
- (c) Finally, we develop an approach for *interactively compiling* high-level conceptual specifications into statistical models.

1.3 Thesis outline

Chapter 2 covers related work that contextualizes the above contributions. The remainder of the dissertation describes how through iterative system development and empirical studies, we came to develop new domain-specific languages (DSLs), representations, and reasoning approaches for authoring statistical analyses.

Chapter 3 presents Tea, a DSL and runtime system for Null Hypothesis Significance Testing (NHST). After discussing more specific related work and explaining the rationale behind supporting NHST (Section 3.1), the chapter describes a usage scenario that illustrates how an analyst would use Tea and how it differs from existing tools (Section 3.2), discusses key design considerations to improve statistical testing practice (Section 3.3), describes the DSL (Subsection 3.4.1) and constraint-based runtime system (Subsection 3.4.2), evaluates Tea against a corpus of expert test choices and a naive test selection regime (Section 3.5), and briefly discusses the limitations and opportunities for future work (Section 3.6). The chapter concludes with a summary of how our work on Tea furthers the thesis of this dissertation.

Chapter 4 introduces our theory of *hypothesis formalization*. While Tea established the feasibility and benefits of designing a DSL focused on capturing implicit data assumptions and hypotheses and developing a formal model of statistical test selection, this chapter steps back to describe data analysis more holistically. This chapter retrospectively justifies our design in Tea and directly informs our work on Tisane, the following chapter. Chapter 4 connects hypothesis formalization to characterizations of “statistical thinking” and situates data analysis in the larger context of scientific discovery (Section 4.1). The chapter proceeds to describe a content analysis that sensitized us to key hypothesis formalization steps (Section 4.2), a lab study observing data analysts *in situ* (Section 4.3), and a qualitative assessment of existing statistical analysis tools (Section 4.4). Based on

these empirical studies, we derive three design implications for how tools can facilitate hypothesis formalization (Section 4.5) and discuss what problem solving strategies (and shortcuts) analysts employ without explicit support for hypothesis formalization (Section 4.6). This chapter also concludes with a summary of how the theory of *hypothesis formalization* informs the thesis.

Chapter 5 is the best representation of how this dissertation grapples with an understanding of data analysis practices (i.e., *hypothesis formalization*), statistical methods, and empirical evidence for what analysts find usable in order to iteratively design and evaluate a DSL and interactive disambiguation process. After covering related work and background on linear modeling (Section 5.1), this chapter describes the first version of the Tisane DSL (Section 5.3), case studies evaluating Tisane (Section 5.4), a study to refine the DSL (Section 5.6), the second major iteration released as rTisane (Section 5.7), and a controlled lab study (Section 5.8). The controlled lab study serves as a summative evaluation of the key tenets of this dissertation. The chapter concludes with key insights derived from iteratively designing and evaluating Tisane and rTisane as well as a few immediate next steps for improving rTisane based on study findings.

Finally, Chapter 6 revisits the key challenges of the thesis and how the projects in this dissertation address each (Section 6.1). Chapter 6 also briefly discusses the real-world impact the DSLs developed in this dissertation have had, offering another form of evidence in support of the thesis. The chapter also discusses how the projects in this dissertation and recent developments in AI-driven code generation create new opportunities to make data analysis authoring valid-by-design and more approachable for statistical non-experts (Section 6.3).

Chapter 2

Related work

This dissertation builds on theories of sensemaking, empirical findings on current analytical praxis, and existing tools throughout the data lifecycle. Additionally, this dissertation uses Donald Campbell’s theory of validity to motivate system designs and interpret evaluation results. Subsequent sections provide additional background as applicable.

2.1 Statistical data analysis as sensemaking

Human beings engage in *sensemaking* to acquire new knowledge. Several theories of sensemaking [PC05; RSPC93; KPRP07] describe how and when human beings seek and integrate new data (e.g., observations, experiences, etc.) to develop their mental models about the world.

Russell et al. [RSPC93] define sensemaking as “the process of searching for a representation and encoding data in that representation to answer task-specific questions.” Russell et al. emphasize the importance of external representations. Sensemaking is the iterative process of searching for and refining external representations in a “learning loop complex” that involves transitioning back and forth between (i) searching for and (ii) instantiating representations. External representations are critical because they influence the quality of conclusions reached at the end of the sensemaking process and affect how much time and effort is required in the process. Some representations may lead to insights more quickly. Indeed, we posit and find that statistical analysis, specifically hypothesis formalization (Chapter 4), is a learning loop [RSPC93] where the conceptual research question or hypothesis is an external representation of a set of assumptions analysts may have about the world (e.g., an implicit causal model), that ultimately affects which statistical models are implemented and which results are obtained. We also find that there are smaller learning loops—for revising explicit causal models, mathematical equations, and partially specified models—embedded in the larger loop of hypothesis formalization.

Grolemund and Wickham argued for statistical data analysis as a sensemaking activity [GW14]. They emphasize the (1) bidirectional nature of updating mental models of the world and hypotheses based on data and collecting data based on hypotheses and (2) the process of identifying and reconciling discrepancies between hypotheses and data. Similar to Russell et al., Grolemund and Wickham’s model demonstrates the importance of representing and re-representing conceptual knowledge. Grolemund and Wickham’s theory of data analysis includes a back and forth between an analyst’s “schema” of how a phenomenon occurs in the world, a statistical model, and data. Analysts’ domain expertise influence their schemas, which represent conceptual knowledge about known and unknown causal mechanisms, for example. Analysts’ conceptual schema directly inform their hypotheses, which are statistical predictions represented in statistical models. These statistical models are then compared to collected data, and any discrepancies between the data and hypothesis require analysts to re-examine and possibly update their statistical model, schema, or both. Extending Grolemund and Wickham’s model, our work on hypothesis formalization differentiates between conceptual and statistical hypotheses and probes the phases an analyst must go through to encode a conceptual hypothesis into a statistical model.

Given the centrality of external representations of implicit conceptual knowledge to authoring statistical analyses that help analysts make sense of the world, we argue that our statistical software should focus on helping analysts to express their conceptual hypotheses and implicit domain knowledge. Through the development of two software systems, Tea (Chapter 3) and Tisane(Chapter 5), we explore *how* to design programming abstractions and *what* those abstractions should include in order for statistical non-experts to externalize their implicit conceptual knowledge about a domain.

2.2 Empirical accounts of data analysis practice

Data analysis involves a number of tasks that involve data discovery, wrangling, profiling, modeling, and reporting [KPHH12]. Extending the findings of Kandel et al. [KPHH12], both Alspaugh et al. [AZL⁺18] and Wongsuphasawat et al. [WLH19] propose exploration as a distinct task. Whereas Wongsuphasawat et al. argue that exploration should subsume discovery and profiling, Alspaugh et al. describe exploration as an alternative to modeling. The importance of exploration and its role in updating analysts’ understanding of the data and their goals and hypotheses is of note, regardless of the precise order or set of tasks. Battle and Heer describe exploratory visual analysis (EVA), a subset of exploratory data analysis (EDA) where visualizations are the primary interfaces and outputs for exploring data, as encompassing both data-focused (bottom-up) and goal- or hypothesis-focused (top-down) investigations [BH19]. In Chapter 4, we found that (i) analysts explored their data before modeling and (ii) exploratory observations sometimes prompted conceptual shifts in

hypotheses (bottom-up) but at other times were guided by hypotheses and only impacted statistical analyses (top-down). In this way, data exploration appears to be an important intermediate step in hypothesis formalization, blurring the lines between exploratory and confirmatory data analysis.

Decisions throughout analysis tasks can give rise to a “garden of forking paths” [GL13], which compounds for meta-analyses synthesizing previous findings [KKH19]. Liu, Boukhelifa, and Eagan [LBE19] proposed a broad framework that characterizes analysis alternatives using three different *levels of abstraction*: cognitive (e.g., shifts in conceptual hypotheses), artifact (e.g., choice in statistical tools), and execution (e.g., computational tuning). *Cognitive* alternatives involve more conceptual shifts and changes (e.g., mental models, hypotheses). *Artifact* alternatives pertain to tooling (e.g., which software is used for analysis?), model (e.g., what is the general mathematical approach?), and data choices (e.g., which dataset is used?). *Execution* alternatives are closely related to artifact alternatives but are more fine-grained programmatic decisions (e.g., hyperparameter tuning). We find that hypothesis formalization involves all three levels of abstraction and provide a more granular depiction of how these levels cooperate with one another (Chapter 4).

Moreover, Liu, Althoff, and Heer [LAH19] identified numerous decision points throughout the data lifecycle, which they call *end-to-end analysis*. They found that analysts often revisit key decisions during data collection, wrangling, modeling, and evaluation. Liu, Althoff, and Heer also found that researchers executed and selectively reported analyses that were already found in prior work and familiar to the research community. The focus of this thesis is on how any single pass or iteration occurs. We approach this work from the perspective that by understanding a single iteration, we may be able to focus analysts on their iterations that are most substantial and impactful and eliminate a number of unnecessary iterations that arise due to mistakes in aligning conceptual and statistical concerns, which we found in our case studies (see Section 5.4).

Importantly, our work differs in (i) scope and (ii) method from prior work in HCI on data analysis practices. Whereas translating a research question or hypothesis into a statistical analysis has remained implicit in prior descriptions of data analysis, we explicate this specific process. Additionally, while previous researchers have relied primarily on post-analysis interviews with analysts, our lab study (Section 4.3) enables us to observe decision making during this process in-situ.

2.3 Tools for data analysis

The software ecosystem for data analysis is vibrant, with numerous programming languages, software packages, and graphical-first tools. A common limitation of existing software is its siloing of statistical specification from the conceptual and data collection details that inadvertently influence statistical analysis. In contrast, the systems in this dissertation explore ways to leverage implicit

conceptual and data collection knowledge to derive statistical analyses. Below, we compare and contrast this dissertation with existing software for conceptual modeling, study design, and statistical specification.

2.3.1 Tools for conceptual modeling

For statistical experts, causal diagramming is a common approach to externalizing implicit conceptual models. For instance, Daggity [THK11] supports authoring, editing, and formally analyzing causal graphs through code and a visual editor. The key limitation of Daggity is that it requires analysts to specify a formal causal graph, which statistical non-experts, including many domain experts, may not be able to do [SSY20; SV18; VDN⁺13]. In fact, an open challenge for causal reasoning and discovery is in getting domain experts to express their implicit knowledge in a way that can be formally represented and reasoned about. Our work on Tisane directly addresses this challenge. Moreover, even if analysts are able to express causal diagrams in Dagitty, Dagitty does not translate queries analysts may have about the causal diagram (i.e., research questions, hypotheses) into statistical models that could assess specific relationships of interest. Tisane also overcomes this limitation for a set of queries about average causal effect.

2.3.2 Tools for study design

Several domain-specific languages [SH17; BEB14], software packages [Tan21; BCCH19a], and standalone applications [MABL⁺07; EWBLM19] specialize in experiment design. A primary focus is to provide researchers low-level control over trial-level and randomization details. For example, JsPsych [DL15a] gives researchers fine-grained control over the design and presentation of stimuli for online experiments. At a mid-level of abstraction, Touchstone [MABL⁺07] is a tool for designing and launching online experiments. It also refers users to R and JMP for data analysis but does not help users author an appropriate statistical model. Touchstone2 [EWBLM19] helps researchers design experiments based on statistical power. At a high-level of abstraction, `edibble` [Tan21] helps researchers plan their data collection schema. `edibble` aims to provide a “grammar of study design” that focuses users on their experimental manipulations in relation to specific units (e.g., participants, students, schools), the frequency and distribution of conditions (e.g., within-subjects vs. between-subjects), and measures to collect (e.g., age, grade, location) in order to output a table to fill in during data collection. While Tisane’s study design specification language uses an abstraction level comparable to `edibble`, Tisane is focused on using the expressed data measurement relationships to infer a statistical model. Additionally, Tisane’s SDSL provides conceptual relationships that are out of the scope of `edibble` but important for specifying conceptually valid statistical models.

2.3.3 Tools for statistical specification

A contribution of this thesis is a closer examination of how existing statistical analysis tools fail to support authoring (Section 4.4). Here, we contrast the systems developed in this thesis to discipline-specific software tools for research and more general automated statistics approaches.

Research has introduced tools to support statistical analysis in diverse domains. ExperiScope [GDH07] supports users in analyzing complex data logs for interaction techniques. ExperiScope surfaces patterns in the data that would be difficult to detect manually and enables researchers to collect noisier data in the wild that have greater external validity. Statsplorer [WSVB15] is an educational web application for novices learning about statistics. While more focused on visualizing various alternatives for statistical tests, Statsplorer automates test selection (for a limited number of statistical tests and by executing simple switch statements) and the checking of assumptions (though it is currently limited to tests of normality and equal variance). [WSVB15] found that Statsplorer helps HCI students perform better in a subsequent statistics lecture. Similar in scope to Statsplorer, Tea is designed to help statistical non-experts author Null-Hypothesis Significance Tests. Tea supports twice as many statistical tests as Statsplorer, suggesting that Tea’s constraint-based approach is more expressive than Statsplorer’s decision-tree implementation for statistical test selection. In contrast to the above systems, a key design consideration for Tea and Tisane has been their ability to apply widely across disciplines and integrate into many existing workflows. Therefore, the systems are implemented as embedded DSLs in Python and R, two widely used programming languages for data science.

The Automatic Statistician [LDG⁺14] generates a report listing all “interesting” relationships (e.g., correlations, statistical models, etc.). Although apparently complete, the Automatic Statistician may overlook analyses that are conceptually interesting and difficult, if not impossible, to deduce from data alone. Furthermore, AutoML tools such as Auto-WEKA [THHLB13], auto-sklearn [FKE⁺15], and H2O AutoML [LP20] also prioritize finding patterns in data and aim to make statistical methods more widely usable. However, Tea and Tisane differ from AutoML efforts in their researchers developing scientific theories. As a result, Tisane provides focus on analysts who prioritize explanation, not just prediction, such as support for specifying GLMMs, which some prominent AutoML tools, such as auto-sklearn [FKE⁺15], omit.

2.4 Validity in statistical data analysis

Finally, a aspect of this thesis is that software with conceptually grounded programming abstractions and automated reasoning can improve the validity of analyses. There are many working definitions of “validity,” from predictive accuracy to a quality of how well experiments are designed

to a trade-off between model simplicity and fit (e.g., R-squared). Donald Campbell’s theory of validity [Sha10], widely adopted across disciplines, provides a framework for reasoning about and unifying many intuitive definitions of validity. Campbell defines four dimensions of validity: internal validity, external validity, statistical conclusion validity, and construct validity. This thesis focuses on enhancing statistical conclusion, external, and internal validity through the correct application and specification of statistical analyses that match analysts’ intentions (i.e., their research questions and hypotheses) and data collection procedures. We do not address construct validity because construct validity is specific to a discipline’s theories and is often debated over a relatively long period of time. In the conclusion (Chapter 6), we touch upon opportunities for future work to address construct validity through the application of recent natural language processing advances.

Chapter 3

Tea: A Domain-Specific Language and Runtime System for Hypothesis Testing

The enormous variety of modern quantitative methods leaves researchers with the non-trivial task of matching analysis and design to the research question.

- Ronald Fisher [Fis37]

Since the development of modern statistical methods (e.g., Student's t-test, ANOVA, etc.), statisticians have acknowledged the difficulty of identifying which statistical tests people should use to answer their specific research questions. Almost a century later, choosing appropriate statistical tests for evaluating a hypothesis remains a challenge. As a consequence, errors in statistical analyses are common [KR12], especially given that data analysis has become a common task for people with little to no statistical expertise.

A wide variety of tools (such as SPSS [Wik19d], SAS [Wik19c], and JMP [Wik19a]), programming languages (e.g., R [Wik19b]), and libraries (including numpy [Oli06], scipy [JOP⁺21a], and statsmodels [SP10]), enable people to perform specific statistical tests, but they do not address the fundamental problem that users may not know which statistical test to perform and how to verify that specific assumptions about their data hold. In fact, all of these tools place the burden of valid, replicable statistical analyses on the user and demand deep knowledge of statistics.

Users not only have to identify their research questions, hypotheses, and domain assumptions, but also must select statistical tests for their hypotheses (e.g., Student's t-test or one-way ANOVA). For each statistical test, users must be aware of the statistical assumptions each test makes about the data (e.g., normality or equal variance between groups) and how to check for them, which requires additional statistical tests (e.g., Levene's test for equal variance), which themselves may demand further assumptions about the data. This cognitively demanding process requires significant

knowledge about statistical tests and their preconditions as well as the ability to perform the tests and verify their preconditions. This process can easily lead to mistakes.

In response, we design and developed Tea¹, a high-level declarative language for automating statistical test selection and execution that abstracts the details of statistical analysis from the users. Tea captures users' hypotheses and domain knowledge, translates this information into a constraint satisfaction problem, identifies all valid statistical tests to evaluate a hypothesis, and executes the tests. Tea's higher-level, declarative nature aims to lower the barrier to valid, replicable analyses.

Tea is easy to integrate directly into common data analysis workflows for users who have minimal programming experience. Tea is implemented as an open-source Python library, so programmers can use Tea wherever they use Python, including within Python notebooks.

In addition, Tea is flexible. Its abstraction of the analysis process and use of a constraint solver to select tests is designed to support its extension to emerging statistical methods, such as Bayesian analysis. Currently, Tea supports frequentist Null Hypothesis Significance Testing (NHST).

This work makes the following contributions:

- a novel DSL for automatically selecting and executing statistical analyses based on users' hypotheses and domain knowledge (Subsection 3.4.1),
- a runtime system that formulates statistical test selection as a maximum constraint satisfaction problem (Subsection 3.4.2), and
- an initial evaluation showing that Tea can express and execute common NHST statistical tests (Section 3.5).

After describing related work, the chapter describes a usage scenario, providing an overview of Tea (Section 3.2). Then, we discuss the concerns about statistics in the HCI community that shaped Tea's design (Section 3.3), the implementation of Tea's programming language (Subsection 3.4.1), the implementation of Tea's runtime system (Subsection 3.4.2), and the evaluation of Tea as a whole (Section 3.5). The chapter concludes with a discussion of Tea's goals, limitations, and future work (Section 3.6) and a summary of how Tea demonstrates my thesis(Section 3.7)

3.1 Background and Related work

Domain-specific languages encapsulate key, routine ideas of domain (e.g., statistical analysis), making programs more concise to write for end-users, providing interfaces to connect with other DSLs

¹named after Fisher's "Lady Tasting Tea" experiment [Fis37]

and systems, and shift the burden of accurate processing from users to systems through specialized reasoning. In the context of the data lifecycle, researchers have developed DSLs that focus on supporting various stages of data exploration, experiment design, and data cleaning. To support data exploration, Vega-lite [SMWH17] is a high-level declarative language that supports users in developing interactive data visualizations without writing functional reactive components. PlanOut [BEB14] is a DSL for expressing and coordinating online field experiments. More niche than PlanOut, Touchstone2 provides the Touchstone Language for specifying condition randomization in experiments (e.g., Latin Squares) [EWBLM19]. To support rapid data cleaning, Wrangler [KPHH11] combines a mixed-initiative interface with a declarative transformation language. Tea provides a language to support another crucial step in the data life cycle: statistical analysis. Tea can be integrated into data analysis workflows and work in tandem with tools such as Wrangler that produce cleaned CSV files ready for analysis.

Furthermore, languages provide semantic structure and meaning that can be reasoned about automatically. For domains with well defined goals, constraint solvers can be a promising technique. Some of the previous constraint-based systems in HCI have been Scout [SKF18], a mixed-initiative system for rapidly prototyping interface designs. Designers specify high-level constraints based on design concepts (e.g., a profile picture should be more emphasized than the name), and Scout synthesizes novel interfaces. Scout also uses Z3’s theories of booleans and integer linear arithmetic. More specific to the data lifecycle are Draco [MWN⁺19] and SetCoLa [HBH18], which formalize visualization constraints for graphs. Whereas SetCoLa is specifically focused on graph layout, Draco formalizes visualization best practices as logical constraints to synthesize new visualizations. The knowledge base can grow and support new design recommendations with additional constraints. Similarly, Tea codifies tests and their preconditions as constraints in a knowledge base. Tea aims to provide an architecture that supports the growth of a statistical analysis knowledge base as communities adopt new statistical best practices and methods. To our knowledge, Tea is the first constraint-based system for statistical analysis.

Statistical scope

Tea is designed for statistical tests common to Null Hypothesis Significance Testing (NHST). While there are calls to incorporate other methods of statistical analysis [KNH16; KR12], Null Hypothesis Significance Testing (NHST) remains the norm in HCI and other disciplines. Therefore, Tea currently implements a module for NHST with the tests found to be most common by [WSVB15]. In particular, Tea supports four classes of tests: correlation (parametric: Pearson’s r , Pointbiserial; non-parametric: Kendall’s τ , Spearman’s ρ), bivariate mean comparison (parametric: Student’s

```

import tea
tea.data('USCrime.csv') 1
variables = [
    {
        'name' : 'So',
        'data type' : 'nominal',
        'categories' : ['0', '1']
    },
    {
        'name' : 'Prob',
        'data type' : 'ratio',
        'range' : [0,1]
    }
]
tea.define_variables(variables) 2
study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design) 3
assumptions = {
    'groups normally distributed': [['So', 'Prob']],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions) 4
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis) 5

```

Figure 3.1: Sample Tea program. The specification outlines an experiment to analyze the relationship between geographic location ('So') and probability of imprisonment ('Prob') in a common USCrime data set [VR13; Kab11]. See Section 3.2 for an explanation of the code. Tea programs specify 1) data, 2) variables, 3) study design, 4) assumptions, and 5) hypotheses.

t-test, Paired t-test; non-parametric: Mann-Whitney U, Wilcoxon signed rank, Welch's), multivariate mean comparison (parametric: F-test, Repeated measures one way ANOVA, Factorial ANOVA, Two-way ANOVA; non-parametric: Kruskal Wallis, Friedman), and comparison of proportions (Chi Square, Fisher's Exact). Tea also supports an implementation of bootstrapping [Efr92].

3.2 Usage Scenario

This section describes how an analyst can use Tea to answer their research questions. We use as an example analyst a historical criminologist who wants to determine how imprisonment differed across regions of the US in 1960². Figure 3.1 shows the Tea code for this example.

The analyst specifies the data file's path in Tea. Tea handles loading and storing the data set for the duration of the analysis session. The analyst does not have to worry about reformatting the data during the analysis process in any way.

The analyst asks if the probability of imprisonment was higher in southern states than in non-southern states. The analyst identifies two variables that could help them answer this question: the probability of imprisonment ('Prob') and geographic location ('So'). Using Tea, the analyst defines

²The example is taken from Ehrlich [Ehr73] and Vandaele [Van87]. The data set comes as part of the MASS package in R.

the geographic location as a dichotomous nominal variable where ‘1’ indicates a southern state and ‘0’ indicates a non-southern state, and indicates that the probability of imprisonment is a numeric data type (ratio) with a range between 0 and 1.

The analyst then specifies their study design, defining the study type to be “observational study” (rather than “experimental study”) and defining the contributor (independent) variable to be the geographic location and the outcome (dependent) variable to be the probability of imprisonment.

Based on their prior research, the analyst knows that the probability of imprisonment in southern and non-southern states is normally distributed. The analyst provides an assumptions clause to Tea in which they specify this domain knowledge. They also specify an acceptable Type I error rate (probability of finding a false positive result), more colloquially known as the ‘significance threshold’ ($\alpha = .05$) that is acceptable in criminology. If the analyst does not have assumptions or forgets to provide assumptions, Tea will use the default of $\alpha = .05$.

The analyst hypothesizes that southern states will have a higher probability of imprisonment than non-southern states. The analyst directly expresses this hypothesis in Tea. *Note that at no point does the analyst indicate which statistical tests should be performed.*

From this point on, Tea operates entirely automatically. When the analyst runs their Tea program, Tea checks properties of the data and finds that the Student’s t-test is appropriate. Tea executes the Student’s t-test and non-parametric alternatives, such as the Mann-Whitney U test, which provide alternative, consistent results.

Tea generates a table of results from executing the tests, ordered by their power (i.e., results from the parametric t-test will be listed first given that it has higher power than the non-parametric equivalent). Based on this output, the analyst concludes that their hypothesis—that the probability of imprisonment was higher in southern states than in non-southern states in 1960—is supported. The results from alternative statistical tests support this conclusion, so the analyst can be confident in their assessment.

The analyst can now share their Tea program with colleagues. Other researchers can easily see what assumptions the analyst made and what the intended hypothesis was (since these are explicitly stated in the Tea program), and reproduce the exact results using Tea.

3.3 Design Considerations

In designing Tea’s language and runtime system, we considered best practices for conducting statistical analyses and derived our own insights on improving the interaction between users and statistical tools.

We identified five key recommendations for statistical analysis from Cairns’ report on common

Table 3.1: Comparison of Tea to other tools. Despite the published best practices for statistical analyses, most tools do not help users select appropriate tests. Tea not only addresses the best practices but also supports reproducing analyses.

| Best practices | SAS | SPSS | JMP | R | Statsplorer [WSVB15] |
|--|-----|-----------|-----------|-----------|----------------------|
| Explicit statement of user assumptions | — | — | — | — | ✓ |
| Automatic verification of test preconditions | — | — | sometimes | sometimes | ✓ |
| Automatic accounting of multiple comparisons | — | — | — | — | — |
| Surface alternative analyses | — | — | — | — | ✓ |
| Contextualize results | ✓ | sometimes | ✓ | sometimes | ✓ |
| Easy to reproduce analysis | ✓ | ✓ | — | ✓ | — |

statistical errors in HCI [Cai07], which echoes many concerns articulated by Wilkinson [Wil99], and from the American Psychological Association’s Task Force on Statistical Inference [Ass96]:

- Users should make explicit their assumptions about the data [Ass96].
- Users should verify and report the results from checking assumptions statistical tests make about the data and variables [Cai07; Ass96].
- Users should account for multiple comparisons [Cai07; Ass96].
- When possible, users should consider alternative analyses that test their hypothesis and select the simplest one [Ass96].
- Users should contextualize results from statistical tests using effect sizes and confidence intervals [Ass96].

An additional practice we wanted to simplify in Tea was *reproducing analyses*. Table 3.1 shows how Tea compares to current tools in supporting these best practices.

Based on these guidelines, we identified two key interaction principles for Tea:

1. *Users should be able to express their expertise, assumptions, and intentions for analysis.* Users have domain knowledge and goals that cannot be expressed with the low-level API calls to the specific statistical tests required by the majority of current tools. A higher level of abstraction that focuses on the goals and context of analysis is likely to appeal to users who may not have statistical expertise (Subsection 3.4.1).
2. *Users should not be burdened with statistical details to conduct valid analyses.* Currently, users must not only remember their hypotheses but also identify possibly appropriate tests and manually check the preconditions for all the tests. Simplifying the user’s procedure by automating the test selection process can help reduce cognitive demand (Subsection 3.4.2).

While there are calls to incorporate other methods of statistical analysis [KNH16; KR12], Null Hypothesis Significance Testing (NHST) remains the norm in HCI and other disciplines. Therefore, Tea currently implements a module for NHST with the tests found to be most common by [WSVB15].

3.4 System overview

Tea consists of a high-level DSL and a runtime system. There are three key steps to compiling a Tea program from user specifications to executing statistical tests:

1. **Check for completeness and syntax.** Tea first checks that a user’s program specifies a data set, variable declarations, study design description, a set of assumptions, and hypotheses using the correct syntax. The data set can be empty (with only column names), which may be useful for pre-registration for instance. If there are any syntax errors or missing parts, Tea will issue an error and stop execution.
2. **Check for consistent, well-formed hypotheses.** Using the variable declarations, Tea then checks that the hypotheses the user states are consistent with variable data types. For instance, Tea would issue an error and halt execution if a nominal variable was hypothesized to have a positive relationship with another nominal variable. If the nominal variables have categories given by numbers (e.g., a variable for education where ‘1’ stands for ‘High School’, ‘2’ for ‘College’, etc.), a linear relationship would be possible to compute by treating the categories as raw continuous values. However, treating the numbers as values is incorrect and the results misleading because the numbers represent discrete categories, not continuous values. Tea avoids such mistakes.
3. **Inspect data properties and infer valid statistical tests.** Once Tea’s compiler verifies that a Tea program is complete, syntactically correct, and consistent, Tea’s runtime system inspects the data to verify properties about it and find a set of valid statistical tests. The higher-level Tea program is then compiled to logical constraints, which is further discussed in Subsection 3.4.2.

3.4.1 Tea’s Domain-Specific Language

Tea is a DSL embedded in Python, implemented as a Python library³. It takes advantage of existing Python data structures (e.g., classes, dictionaries, and enums). We chose Python because of its widespread adoption in data science.

³Tea is open-source and available for download on pip, a common Python package manager.

A key challenge in designing Tea’s DSL is determining the level of granularity necessary to produce an accurate analysis. In Tea programs, users describe their studies in five ways: (1) providing a data set, (2) describing the variables of interest in that data set, (3) specifying their study design, (4) stating their assumptions about the variables, and (5) formulating hypotheses about the relationships between variables. Figure 3.2 shows an example Tea program and its output.

Data

Data is required for executing statistical analyses. One challenge in managing data for analysis is minimizing both duplicated data and user intervention.

To reduce the need for user intervention for data manipulation, Tea requires the data to be a CSV in long format. CSVs are a common output format for data storage and cleaning tools. Long format (sometimes called “tidy data” [W⁺14]) is a denormalized format that is widely used for collecting and storing data, especially for within-subjects studies.

Unlike R and Python libraries such as numpy [Oli06], Tea only requires one instance of the data. Users do not have to duplicate the data or subsets of it for analyses that require the data to be in slightly different forms. Minimizing data duplication or segmentation is also important to avoid user confusion about where some data exist or which subsets of data pertain to specific statistical tests.

Optionally, users can also indicate a column in the data set that acts as a relational (or primary) key, or an attribute that uniquely identifies rows of data. For example, this key could be a participant identification number in a behavioral experiment. A key is useful for verifying a study design, described below. Without a key, Tea’s default is that all rows in the data set comprise independent observations (that is, all variables are between subjects).

To use Tea for pre-registration prior to collecting data, a CSV with only column names is necessary.

Variables

Variables represent columns of interest in the data set. Variables have a name, a data type (*nominal*, *ordinal*, *interval*, or *ratio*), and, when appropriate, valid categories. Users (naturally) refer to variables through a Tea program using their names. Only nominal and ordinal variables have a list of possible categories. For ordinal variables, the categories are also ordered from left to right.

Variables encapsulate queries. The queries represent the index of the variable’s column in the original data set and any filtering operations applied to the variable. For instance, it is common to filter by category for nominal variables.

Study Design

Three aspects of study design are important for conducting statistical analyses: (1) the type of study (observational study vs. randomized experiment), (2) the independent and dependent variables, and (3) the number of observations per participant (e.g., between-subjects variables vs. within-subjects variables).

For semantic precision, Tea uses different terms for independent and dependent variables for observational studies and experiments. In experiments, variables are described as either “independent” or “dependent” variables. In observational studies, variables are either “contributor” (independent) or “outcome” (dependent) variables.

Assumptions

Users’ assumptions based on domain knowledge are critical for conducting and contextualizing studies and analyses. Often, users’ assumptions are particular to variables and specific properties (e.g., equal variances across different groups). Current tools generally do not require that users encode these assumptions, leaving them implicit.

Tea takes the opposite approach to contextualize and increase the transparency of analyses. It requires that users be explicit about assumptions and statistical properties pertaining to the analysis as a whole (e.g., acceptable Type I error rate/significance threshold) and the data.

Tea supports two modes for treating user assumptions: *strict* and *relaxed*. In both modes, Tea verifies all user assumptions and issues warnings for assumptions that statistical testing does not verify. In the *strict* mode, Tea overrides user assumptions when selecting valid statistical tests. In the *relaxed* mode, Tea defers to user assumptions and proceeds as if the assumptions verified even if they did not. The *strict* mode is the default, but users can specify the *relaxed* mode. Figure 3.2 shows the two modes and the different warnings and output they generate.

If users also know that a data transformation (i.e., log transformation) applies to a variable, they can express this as an assumption. Data transformations are not properties to be verified but rather treatments of data that are applied during assumption verification, statistical test selection, and test execution, which is why they are included in the assumptions clause. The next section discusses the verification process for assumptions in greater detail.

Hypotheses

Hypotheses drive the statistical analysis process. Users often have hypotheses that are technically alternative hypotheses.

Tea focuses on capturing users’ alternative hypotheses about the relationship between two or

more variables. Tea uses the alternate hypothesis to conduct either a two-sided or one-sided statistical test. By default, Tea uses the null hypothesis that there is no relationship between variables.

3.4.2 Tea’s Constraint-based Runtime System

Tea compiles programs into logical constraints about the data and variables, which it resolves using a constraint solver. A significant benefit of using a constraint solver is extensibility. Adding new statistical tests does not require modifying the core of Tea’s runtime system. Instead, defining a new test requires expressing a single new logical relationship between a test and its preconditions.

At runtime, Tea invokes a solver that operates on the logical constraints it computes to produce a list of valid statistical tests to conduct. This process presents three key technical challenges: (1) incorporating statistical knowledge as constraints, (2) expressing user assumptions as constraints, and (3) recursively selecting statistical tests to verify preconditions of other statistical tests.

SMT Solver

As its constraint solver, Tea uses Z3 [DMB08], a Satisfiability Modulo Theory (SMT) solver.

Satisfiability is the process of finding an assignment to variables that makes a logical formula true. For example, given the logical rules $0 < x < 100$ and $y < x$, $\{x = 1, y = 0\}$, $\{x = 10, y = 5\}$, and $\{x = 99, y = -100\}$ would all be valid assignments that satisfy the rules. SMT solvers determine the satisfiability of logical formulas, which can encode boolean, integer, real number, and uninterpreted function constraints over variables. SMT solvers can also be used to encode constraint systems, as we use them here. A wide variety of applications ranging from the synthesis of novel interface designs [SKF18], the verification of website accessibility [PGE⁺18], and the synthesis of data structures [LTE16] employ SMT solvers.

Logical Encodings

The first challenge of framing statistical test selection as a constraint satisfaction problem is defining a logical formulation of statistical knowledge.

Tea encodes the applicability of a statistical test based on its preconditions. A statistical test is applicable if and only if all of its preconditions (which are properties about variables) hold. We derived preconditions for tests from an online HCI and statistics course [KW19], a statistics textbook [FMF12], and publicly available data science resources from universities [Bru19; Lib19].

Tea represents each precondition for a statistical test as an uninterpreted function representing a property over one or more variables. Each property is assigned `true` if the property holds for the variable/s; similarly, if the property does not hold, the property function is assigned `false`.

Tea also encodes statistical knowledge about variable types and properties that are essential to statistical analysis as axioms, such as the constraint that only a continuous variable can be normally distributed.

Algorithm

Tea frames the problem of finding a set of valid statistical tests as a maximum satisfiability (MaxSAT) problem that is seeded with user assumptions.

First, Tea translates each user assumption about a data property into an axiom about a property and variable. As described in Section 3.4.1, user assumptions about properties but not data transformations are always checked. In the *strict* mode, Tea overrides any user assumptions it does not find to hold, creating an axiom that a property is `false`. In the *relaxed* mode, Tea defers to user assumptions, creating axioms that a property is `true`. For any user assumptions that do not pass statistical testing, Tea warns the user and explains how it will proceed depending on the mode.

Then, for each new statistical test Tea tries to satisfy, Tea checks to see if each precondition holds. For each precondition checked, Tea adds the property and variable checked as an axiom to observe as future tests are checked. If any property violates the axioms derived from users' assumptions, the property is removed and the test is invalidated. Users' assumptions always take precedence.

The constraint solver then prunes the search space. Tea does not compute all properties for all variables, a significant optimization when analyzing very large data sets.

At the end of this process, Tea finds a set of valid statistical tests to execute. If this set is empty, Tea defaults to its implementation of bootstrapping [Efr92]. Otherwise, Tea proceeds and executes all valid statistical tests. Tea returns a table of results to users, applying multiple comparison corrections [Hol79] and calculating effect sizes when appropriate.

Optimization: Recursive Queries

When Tea verifies a property holds for a variable, it often must invoke another statistical test. For example, to check that two groups have equal variance, Tea must execute Levene's test. The statistical test used for verification may then itself have a precondition, such as a minimum sample size.

Such recursive queries are inefficient for SMT solvers like Z3 to reason about. To eliminate recursion, Tea lifts some statistical tests to properties. For instance, Tea does not encode the Levene's test as a statistical test. Instead, Tea encodes the property of having equal variance between groups and executes the Levene's test for two groups when verifying that property for particular variables.

User Output

The result of running a Tea program with data is a listing of the results of executing valid statistical tests, as shown in Figure 3.2. For each valid statistical test executed, the output contains the properties of data that Tea checked and used to determine that a statistical test applied, the test statistic value, p-value (and an adjusted p-value, if applicable), effect sizes (Cohen’s d [Coh88] and Vargha Delaney A12 [VD00]), the alpha level the user specified in their program, the precise null hypothesis the statistical test examined, an interpretation of the results in APA format [A⁺83], and text recommending users to focus on effect size rather than the p-value for a holistic view of their data. This output is intended to inform users of why Tea selected specific statistical tests and how to interpret their results.

3.5 Evaluation

We assessed the validity of Tea in two ways. First, we compared Tea’s suggestions of statistical tests to suggestions in textbook tutorials. We use these tutorials as a proxy for expert test selection. Second, for each tutorial, we compared the analysis results of the test(s) suggested by Tea to those of the test suggested in the textbook as well as all other candidate tests. We use the set of all candidate tests as a proxy for non-expert test selection.

We differentiate between *candidate* and *valid* tests. A candidate test can be computed on the data, when ignoring any preconditions regarding the data types or distributions. A valid test is a candidate test for which all preconditions are satisfied.

3.5.1 How does Tea compare to textbook tutorials?

Our goal was to compare Tea to expert recommendations.

We sampled 12 data sets and examples from R tutorials ([Kab11] and [FMF12]). These included eight parametric tests, four non-parametric tests, and one Chi-square test. We chose these tutorials because they appeared in two of the top 20 statistical textbooks on Amazon and had publicly available data sets, which did not require extensive data wrangling.

We translated all analyses into Tea and encoded any assumptions explicitly stated in the tutorial. Tea selected tests based only on the data and the assumptions expressed in the Tea program. Where Tea disagreed with the tutorials, either (1) the tutorial authors chose the wrong analyses or (2) the tutorial authors had implicit assumptions about the data that did not hold up to statistical testing.

For nine out of the 12 tutorials, Tea suggested the same statistical test (see Table 3.2). For three out of 12 tutorials, which used a parametric test, Tea suggested using a non-parametric alternative instead. Tea’s recommendation of using a non-parametric test instead of a parametric one did not

Table 3.2: Results of applying Tea to 12 textbook tutorials.

Tea is comparable to an expert selecting statistical tests. Tea can prevent false positive and false negative results by suggesting only tests that satisfy all assumptions. *Tutorial* gives the test described in the textbook; *Candidate tests (p-value)* gives all tests a user could run on the provided data with corresponding p-values; *Assumptions* gives all satisfied (lightly shaded) and violated (white) assumptions; *Tea suggests* indicates which tests Tea suggests based on their preconditions (assumptions about the data). **Emphasized** p-values indicate instances where a candidate test leads to a wrong conclusion about statistical significance. Although this table focuses on p-values, Tea produces richer output that provides a more holistic view of the statistical analysis results by including effect sizes, for instance. Refer to Figure 3.2 for an example of output from a Tea program.

| Tutorial | | Candidate tests (p-value) | Assumptions* | Tea suggests |
|---|---------------------------------|---------------------------|--------------|--------------|
| Pearson [Kab11] | Pearson's r | (6.96925e-06) | ② ④ ⑤ | — |
| | Kendall's τ | (2.04198e-05) | ② ④ | ✓ |
| | Spearman's ρ | (2.83575e-05) | ② ④ | ✓ |
| Spearman's ρ [FMF12] | Spearman's ρ | (.00172) | ② ④ | ✓ |
| | Pearson's r | (.01115) | ② ④ | — |
| | Kendall's τ | (.00126) | ② ④ | ✓ |
| Kendall's τ [FMF12] | Kendall's τ | (.00126) | ② ④ | ✓ |
| | Pearson's r | (.01115) | ② ④ | — |
| | Spearman's ρ | (.00172) | ② ④ | ✓ |
| Pointbiserial [FMF12] | Pointbiserial (Pearson's r) | (.00287) | ② ④ ⑤ | — |
| | Spearman's ρ | (.00477) | ② ④ | — |
| | Kendall's τ | (.00574) | ② ④ | — |
| | Bootstrap | (<0.05) | | ✓ |
| Student's t-test [Kab11] | Student's t-test | (.00012) | ② ④ ⑤ ⑥ ⑦ ⑧ | ✓ |
| | Mann-Whitney U | (9.27319e-05) | ② ④ ⑦ ⑧ | ✓ |
| | Welch's t-test | (.00065) | ② ④ ⑤ ⑦ ⑧ | ✓ |
| Paired t-test [FMF12] | Paired t-test | (.03098) | ② ④ ⑤ ⑦ ⑧ | ✓ |
| | Student's t-test | (.10684) | ② ④ ⑤ ⑦ | — |
| | Mann-Whitney U | (.06861) | ② ④ ⑦ | — |
| | Wilcoxon signed rank | (.04586) | ② ④ ⑦ ⑧ | ✓ |
| | Welch's t-test | (.10724) | ② ⑦ | — |
| Wilcoxon rank [FMF12] | Wilcoxon signed rank | (.04657) | ② ④ ⑦ ⑧ | ✓ |
| | Student's t-test | (.02690) | ② ④ ⑦ | — |
| | Paired t-test | (.01488) | ② ④ ⑤ ⑦ ⑧ | — |
| | Mann-Whitney U | (.00560) | ② ④ ⑦ | — |
| | Welch's t-test | (.03572) | ② ④ ⑦ | — |
| F-test [FMF12] | F-test | (9.81852e-13) | ② ④ ⑤ ⑥ ⑨ | ✓ |
| | Kruskal Wallis | (2.23813e-07) | ② ④ ⑨ | ✓ |
| | Friedman | (8.66714e-07) | ② ⑦ | — |
| | Factorial ANOVA | (9.81852e-13) | ② ④ ⑤ ⑥ ⑨ | ✓ |
| Kruskal Wallis [FMF12] | Kruskal Wallis | (.03419) | ② ④ ⑨ | ✓ |
| | F-test | (.05578) | ② ④ ⑤ ⑨ | — |
| | Friedman | (3.02610e-08) | ② ⑦ | — |
| | Factorial ANOVA | (.05578) | ② ④ ⑤ ⑨ | — |
| Repeated measures one way ANOVA [FMF12] | Repeated measures one way ANOVA | (.0000) | ② ④ ⑤ ⑥ ⑦ ⑨ | ✓ |
| | Kruskal Wallis | (4.51825e-06) | ② ④ ⑦ ⑨ | — |
| | F-test | (1.24278e-07) | ② ④ ⑤ ⑥ ⑦ ⑨ | — |
| | Friedman | (5.23589e-11) | ② ④ ⑦ ⑨ | ✓ |
| | Factorial ANOVA | (1.24278e-07) | ② ④ ⑤ ⑥ ⑨ | ✓ |
| Two-way ANOVA [FMF12] | Two-way ANOVA | (3.70282e-17) | ② ④ ⑤ ⑨ | — |
| | Bootstrap | (<0.05) | | ✓ |
| Chi Square [FMF12] | Chi Square | (4.76743e-07) | ② ④ ⑨ | ✓ |
| | Fisher's Exact | (4.76743e-07) | ② ④ ⑨ | ✓ |

*① one variable, ② two variables, ③ two or more variables, ④ continuous vs. categorical vs. ordinal data, ⑤ normality, ⑥ equal variance, ⑦ dependent vs. independent observations, ⑧ exactly two groups, ⑨ two or more groups

change the statistical significance of the result at the .05 level. Tea suggested non-parametric tests based on the Shapiro-Wilk test for normality. It is possible that tutorial authors visualized the data to make implicit assumptions about the data, but this practice or conclusion was not made explicit in the tutorials.

For the two-way ANOVA tutorial from [FMF12], which studied how gender and drug usage of individuals affected their perception of attractiveness, a precondition of the two-way ANOVA is that the dependent measure is normally distributed in each category. This precondition was violated. As a result, Tea defaulted to bootstrapping the means for each group and reported the means and confidence intervals. For the pointbiserial correlation tutorial from [FMF12], Tea also defaulted to bootstrap for two reasons. First, the precondition of normality is violated. Second, the data uses a dichotomous (nominal) variable, which invalidates Spearman’s ρ and Kendall’s τ .

Tea generally agrees with expert recommendations and is more conservative in the presence of non-normal data, minimizing the risk of false positive findings.

3.5.2 Does Tea avoid common mistakes made by non-expert users?

Our goal was to assess whether any of the tests suggested by Tea (i.e., valid candidate tests) or any of the invalid candidate tests would lead to a different conclusion than the one drawn in the tutorial. Table 3.2 shows the results. Specifically, emphasized p-values indicate instances for which the result of a test differs from the tutorial in terms of statistical significance at the .05 level.

For all of the 12 tutorials, Tea’s suggested tests led to the same conclusion about statistical significance. For two out of the 12 tutorials, two or more candidate tests led to a different conclusion. These candidate tests were invalid due to violations of independence or normality.

To summarize, the evaluation shows us that (i) Tea can replicate and even improve upon expert choices and (ii) Tea could help novices avoid common mistakes and false conclusions.

3.6 Discussion, Limitations, and Future Work

Our goal with Tea was to determine the feasibility of automating statistical test selection based on high-level input from analysts. Automating statistical test selection raises important concerns about the impact of such automation on the reliability of statistical conclusions. In this regard, there are two chief concerns pertaining to (i) selective inference and (ii) multiple testing, both of which inflate the Type I Error Rate and can lead to more false discoveries.

Tea relies on statistical tests (e.g., Shapiro-Wilk’s test for normality) to assess properties of data to determine which statistical tests (e.g., Student’s t-test) are used to assess the input hypothesis. Repeated property testing of the data is a form of “double-dipping” [], or using the data to make

decisions about analyses on the data. Preventing this would be ideal to reduce the false positive discovery rate. However, the statistics community is still developing techniques to address this issue. A naive approach would be to only use a sample of the data to determine the final statistical test and then use another sample to make statistical inferences. While viable for large datasets, this may not be possible for smaller datasets. A more recently proposed technique, data fission [] overcomes, in theory, this dependence on dataset size. Data fission introduces noise to the data to make analysis decisions (i.e., statistical test selection) and then stripping the noise to obtain final results. Tea does not currently implement either of these approaches. In the future, Tea should incorporate these and future recommendations from the statistics community.

Furthermore, there is an inherent tension between executing multiple statistical tests (e.g., Student's t-test and Welch's t-test) to show analysts the robustness, or sensitivity, of statistical results and increasing the number of comparisons performed. In Tea, we believed that providing analysts with the ability to compare statistical tests, make sensitivity judgments, and report the results of a test most common in their disciplines was more important than restricting the number of statistical tests, especially we have observed analysts intentionally run multiple statistical tests in order to compare results on their own. To more fully support sensitivity analyses and discourage cherry-picking statistical tests and results, Tea should provide more explicit support for interpreting, comparing, and contrasting statistical test results in the future. This will be particularly important in scenarios where statistical tests may disagree with one another. conflicting test conclusions.

Finally, Tea's test selection is well suited for answering a class of relatively simple research questions. At the same time, there are more complex research questions that analysts want to ask about their domain using data that require more complex statistical analyses. These are currently out of reach for Tea. For instance, domain experts may not want to know that there is a difference between treatment and control groups but also estimate the influence of the treatment on an outcome in the presence of other variables that also influence treatment and the outcome. Therefore, in order to support a larger class of research questions and statistical models, we need to re-consider and extent Tea's abstractions and constraint-based reasoning approach.

3.7 Summary of Contributions

A common approach to assessing support for conceptual hypotheses in data is to use statistical tests (e.g., Student's t-test, Chi-Square test, ANOVA). Statistical testing requires analysts to grapple with their conceptual hypotheses, know a number of tests and when they are applicable (i.e., know the preconditions for when these tests hold), assess the applicability of tests (i.e., check preconditions), and pick and implement specific tests using low-level APIs.

Tea’s key insight is that we can reformulate statistical test selection as a constraint satisfaction problem. We designed and implemented a higher-level DSL around this insight that takes an analyst’s hypothesis and assumptions about their data as input and provides the results of executing valid statistical tests as output. In an evaluation, we found that Tea avoids faulty test selection and conclusions that are easy to make using existing tools. In this way, Tea improves statistical conclusion and internal validity [Sha10].

Tea demonstrates the feasibility and benefit of developing systems that emphasize *higher-level abstractions* and *automated reasoning* for statistical tests (Section 1.1). However, using statistics to answer real-world questions requires going beyond statistical testing to grappling with statistical modeling and effect estimation. Next, we consider how our approach generalizes to a larger class of statistical analyses.

This work was done in collaboration with Maureen Daum, Jared Roesch, Sarah E. Chasins, Emery Berger, René Just, and Katharina Reinecke. It was originally published and presented at ACM UIST 2019 cite. Since publication, multiple people, including most notably Shreyash Nigam, Reiden Chea, and Annie Denton, have contributed to updating and improving Tea.

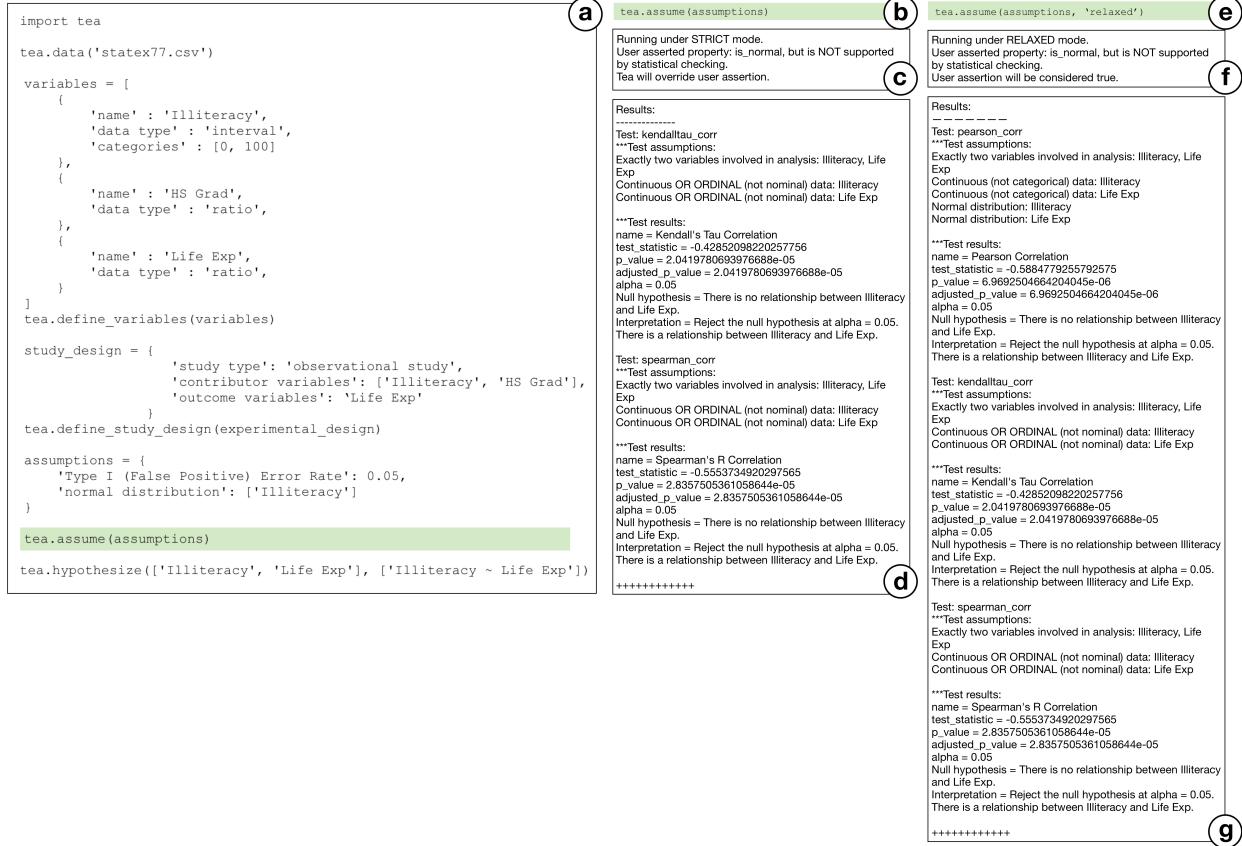


Figure 3.2: Tea program and its mode-dependent executions. a) Tea program that aims to determine if two contributor variables, ‘Illiteracy’ and ‘HS Grad’ that may predict a third outcome variable ‘Life Exp’, are correlated. The user asserts that ‘Illiteracy’ is normally distributed. b) By default, Tea executes programs in the *strict* mode. c) Warning that Tea disagrees with the user and will override the user’s assertion that ‘Illiteracy’ is normally distributed in the *strict* mode. d) Results without the parametric test since Tea overrides user’s assertion. e) A single line change can modify Tea to execute a program in *relaxed* mode. f) Warning that Tea cannot verify normality for ‘Illiteracy’ but will defer to user’s assertion. g) Results with the parametric test since Tea proceeds as if ‘Illiteracy’ was normally distributed.

Chapter 4

Hypothesis Formalization: A conceptual framework describing how analysts translate research questions into statistical analyses

Consider a census researcher who asks, “In the United States (U.S.), how does an individual’s sex relate to their annual income?” Drawing upon their prior experiences and exploratory data visualizations, the researcher knows that income in the U.S. is skewed, and they want to know how the distributions of income among males and females differ (step i). However, before implementing, they (implicitly) define their causal model: The researcher knows that other factors, such as education and race, may be associated with employment opportunities, which may then influence income. As such, they refine their conceptual hypothesis—that sex influences income—to consider the possible effects of race, education, sex, and their interactions on income. They plan to fit a generalized linear model with race, education, sex, and their two-way interactions as predictors of income (step ii). They start implementing a script to load and model data (step iii). The researcher receives a small table of results and is surprised to receive a convergence warning. After further investigation, they simplify their model and remove the interaction effects to see how that may affect convergence (revise step iii). This time, their model’s inference algorithm converges, and they interpret the results (step iv), but they really want to study how sex and race interact, so they return to implementation (step iii) and proceed as before, iteratively removing and adding effects and changing computational parameters, and as a by-product shifting which high-level conceptual hypothesis is reflected in the model.

Performing statistical data analysis goes well beyond invoking the correct statistical functions in a software library. As seen with the census researcher, statistical analyses require (i) translating high-level, domain-specific questions and hypotheses into specific statistical questions [CEG⁺16]; (ii) identifying statistical models to answer the statistical questions; (iii) implementing and executing these statistical models, typically with the help of software tools; and (iv) interpreting the results, considering the domain-specific questions and applying analytical reasoning. Analysts must go back and forth between conceptual hypothesis and model implementation realities, grappling with domain knowledge, limitations of data, and statistical methods.

We refer to the process of translating a conceptual hypothesis into a computable statistical model as *hypothesis formalization*. This process is messy and under-scrutinized in prior work. Consequently, we investigate the steps, considerations, challenges, and tools involved. Based on our findings, we define hypothesis formalization as a dual-search process [KD88] that involves developing and integrating cognitive representations from two different perspectives—conceptual hypotheses and concrete model implementations. Analysts move back and forth between these two perspectives during formalization while balancing conceptual, data-driven, statistical, and implementation constraints. Figure 4.1 summarizes our definition and findings. Specifically, this chapter addresses the following questions to develop our definition of hypothesis formalization:

- **RQ1 - Steps:** What is the range of steps an analyst might consider when formalizing a hypothesis? How do these steps compare to ones that we might expect based on prior work?
- **RQ2 - Process:** How do analysts think about and perform the steps to translate their hypotheses into model implementations? What challenges do they face during this process?
- **RQ3 - Tools:** How might current software tools influence hypothesis formalization?

To sensitive ourselves to the steps (**RQ1 - Steps**) and considerations (**RQ2 - Process**) involved in hypothesis formalization, we compared and contrasted existing models and descriptions of data analysis in prior work. We augmented our deep dive into prior work with a formative content analysis of 50 randomly sampled research papers from five different venues, including Psychological Science and Nature. We find that researchers decompose their research hypotheses into specific sub-hypotheses, derive proxy variables from theory and available data, and adapt statistical analyses to account for data collection procedures. A key takeaway from prior work and the formative content analysis was the “hypothesis refinement loop” in Figure 4.1.

To validate and deepen our understanding of hypothesis formalization (**RQ1 - Steps** and **RQ2 - Process**), we designed and conducted a lab study in which we observed 24 analysts develop and formalize hypotheses in-situ. We find that analysts foreground implementation concerns, even when brainstorming hypotheses, and try to fit their hypotheses and analyses to prior experiences



Figure 4.1: Definition and overview of the hypothesis formalization steps and process. Hypothesis formalization is a dual-search process of translating a **conceptual hypothesis** into a statistical **model implementation**. Blue indicates steps and transitions that we identified. Black indicates steps and transitions discussed in prior work. “Mathematical Equation” (dashed box) was rarely an explicit step in our lab study but evident in our content analysis. Our findings (blue arrows) corroborate and subsume several of the transitions identified in prior work with greater granularity. When they do not, prior work’s transitions are included in black. For example, analysts may operationalize a conceptual hypothesis as a causal model by first decomposing the conceptual hypothesis into sub-hypotheses and then identifying proxy variables to incorporate in a causal model (blue arrows above). Our definition of hypothesis formalization is a consequence of our synthesis of prior work, content analysis, lab study, and analysis of tools. Hypothesis formalization is a non-linear process. Analysts iterate over conceptual steps to refine their hypothesis in a *hypothesis refinement loop*. Analysts also iterate over computational and implementation steps in a *model implementation loop*. Data collection and data properties may also prompt conceptual revisions and influence statistical model implementation. As analysts move toward model implementation, they increasingly rely on software tools, gain specificity, and create intermediate artifacts along the way (e.g., causal models, observations about data, etc.).

and familiar tools, suggesting a strong influence of tools (**RQ3 - Tools**). Thus, the lab study reinforced the hypothesis refinement loop, surfaced the “model implementation loop,” and raised questions about the role of tools.

To identify how tools may shape hypothesis formalization (**RQ3 - Tools**), we reviewed 20 statistical software tools. We find that although the tools support nuanced model implementations, their low-level abstractions can focus analysts on statistical and computational details at the expense of higher-level reasoning about initial hypotheses. Tools also do not aid analysts in identifying reasonable model implementations that would test their conceptual hypotheses, which may explain why analysts in our lab study relied on familiar approaches, even if sub-optimal. Furthermore, our tools review confirmed that the dual processes inform one another during hypothesis formalization.

Taken together, our findings help us define the hypothesis formalization framework, as summarized in Figure 4.1, and suggest **three design implications** for tools to more directly support hypothesis formalization: (i) show the relationships between related statistical models that seem syntactically different from each other, (ii) provide higher-level abstractions for expressing conceptual hypotheses and partial model specifications, and (iii) develop bidirectional computational assistance for authoring causal models and relating them to statistical models.

By defining and characterizing hypothesis formalization, we situate data analysis in a larger model of scientific discovery, identify specific problem solving strategies used in hypothesis formalization that demonstrate how data analysis (and science) is a practice, and identify opportunities for future software to improve the transparency and reproducibility of analyses by explicitly supporting pathways and loops through hypothesis formalization.

4.1 Background and Related Work

Our work integrates and builds up existing theories of statistical thinking in cognitive psychology and statistics. We also situate hypothesis formalization in the larger context of scientific discovery.

4.1.1 Statistical Thinking

Statistical thinking and practice require differentiating between *domain* and *statistical* questions. The American Statistical Association (ASA), a professional body representing statisticians, recommends that universities teach this fundamental principle in introductory courses (see Goal 2 in [CEG⁺16]). Similarly, researchers Wild and Pfannkuch emphasize the importance of differentiating between and integrating statistical knowledge and context (or domain) knowledge when thinking statistically [Pfa97; PW⁺00; WP99]. They propose a four step model for operationalizing ideas (“inklings”) into plans for collecting data, which are eventually statistically analyzed. In their model, analysts must transform “inklings” into broad questions and then into precise questions that are then finally turned into a plan for data collection (see Figure 2 in [WP99]). Statistical and domain knowledge inform all four stages. However, it is unknown what kinds of statistical and domain knowledge are helpful, how they are used and weighed against each other, and when certain kinds of knowledge are helpful to operationalize inkings. Our work in defining hypothesis formalization provides more granular insight into Wild and Pfannkuch’s proposed model of operationalization and aims to answer when, how, and what kinds of statistical and domain knowledge are used during statistical data analysis.

More recently, in *Statistical Rethinking* [McE20], McElreath proposes that there are three key representational phases involved in data analysis: conceptual hypotheses, causal models underly-

ing hypotheses (which McElreath calls “process models”), and statistical models. McElreath, like the ASA and Wild and Pfannkuch, separates domain and statistical ideas and discusses the use of causal models as an intermediate representation to connect the two. McElreath emphasizes that conceptual hypotheses may correspond to multiple causal and statistical models, and that the same statistical model may provide evidence for multiple, even contradictory, causal models and hypotheses. McElreath’s framework does not directly address how analysts navigate these relationships or how computation plays a role, both of which we take up in this chapter.

Overall, our work provides empirical evidence for prior frameworks but also (i) provides more granular insight into *how* and *why* transitions between representations occur and (ii) scrutinizes the role of *software and computation* through close observation of analyst workflows in the lab as well as through a follow-up analysis of statistical software. Based on these observations, we also speculate on how tools might better support hypothesis formalization.

4.1.2 Statistical data analysis as part of scientific discovery

Klahr and Simon characterized scientific discovery as a dual-search process involving the development and evaluation of hypotheses and experiments [KD88]. They posited that scientific discovery involved tasks specific to hypotheses (e.g., revising hypotheses) and to experiments (e.g., analyzing data collected from experiments), which they separated into two different “spaces,” and tasks moving between them, which is where we place hypothesis formalization. Extending Klahr and Simon’s two-space model, Schunn and Klahr proposed a more granular four-space model involving data representation, hypothesis, paradigm, and experiment spaces [SK95; SK96]. In the four-space model, conceptual hypothesizing still lies in the hypothesis space, and hypothesis testing and statistical modeling lies in the paradigm space. As such, hypothesis formalization is a process connecting the hypothesis and paradigm spaces. In Schunn and Klahr’s four-space model, information flows unidirectionally from the hypothesis space to the paradigm space. We extend this prior research with evidence that the path from hypothesis and paradigm spaces is actually bidirectional (see Figure 4.1).

Figure 4.2 augments Schunn and Klahr’s original diagram (Figure 1 in [SK95]) with annotations depicting how our content analysis of research papers and lab study triangulate a tighter dual-space search between hypothesis and paradigm spaces with a focus on hypothesis formalization. Our mixed-methods approach follows the precedent and recommendations of Klahr and Simon’s [KS99] study of scientific discovery activities.

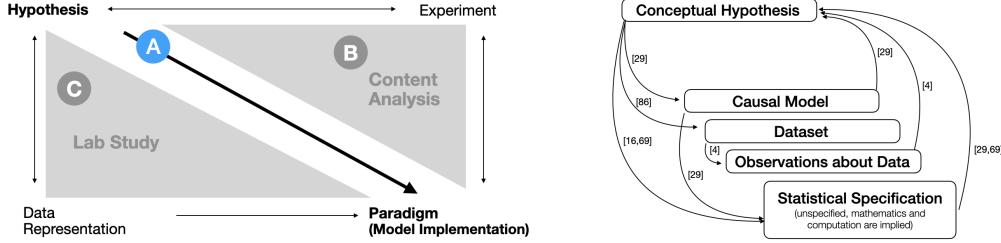


Figure 4.2: Relationship between hypothesis formalization and prior work.

Left: Schunn and Klahr’s four-space model of scientific discovery (stylized adaptation from Figure 1 in [SK95]), which includes unidirectional information flow from the hypothesis space to the paradigm space (which includes model implementation). Hypothesis formalization (A) is focused on a tighter integration and the information flow between hypothesis and paradigm spaces. Specifically, the information flow is bidirectional in hypothesis formalization. Our content analysis (B) and lab study (C) triangulate the four-space model to understand hypothesis formalization from complementary perspectives. *Right:* Hypothesis formalization steps also identified in prior work on theories of sensemaking, statistical thinking, and data analysis workflows (citations included to the right of the arrows). Hypothesis formalization is finer grained and involves more iterations. While prior work broadly refers to mathematical equations, partial model specifications, and computationally tuned model implementations as statistical specifications, hypothesis formalization differentiates them. As a whole, this chapter provides empirical evidence for theorized loops between conceptual hypothesis and statistical specification (see Figure 4.1).

4.2 Formative content analysis

To complement our in-depth synthesis of prior work, we conducted a formative content analysis of 50 peer-reviewed publications from five different domains.

Methods

We randomly sampled ten papers published in 2019 from each of the following venues: (1) the Proceedings of the National Academy of Sciences (PNAS), (2) Nature, (3) Psychological Science (PS), (4) Journal of Financial Economics (JFE), and (5) the ACM Conference on Human Factors in Computing Systems (CHI). We sampled papers that used statistical analyses as either primary or secondary methodologies. Our sample represents a plurality of domains and recent practices.¹

The first two authors iteratively developed a codebook to code papers at the paragraph-level. The codebook contained five broad categories: (i) research goals, (ii) data sample information, (iii) statistical analysis, (iv) results reporting, and (v) computation. Each category had more specific codes to capture more nuanced differences between papers. This tiered coding scheme enabled us to see general content patterns across papers and nuanced steps within papers. The first two authors reached substantial agreement (IRR = .69 - .72) even before resolving disagreements. The first

¹Google Scholar listed the venues among the top three in their respective areas in 2018. Venues were often clustered in the rankings without an obvious top-one, so we chose among the top three based on ease of access to publications (e.g., open access or access through our institution). Some papers were accepted and published before 2019, but the journals had included them in 2019 issues.

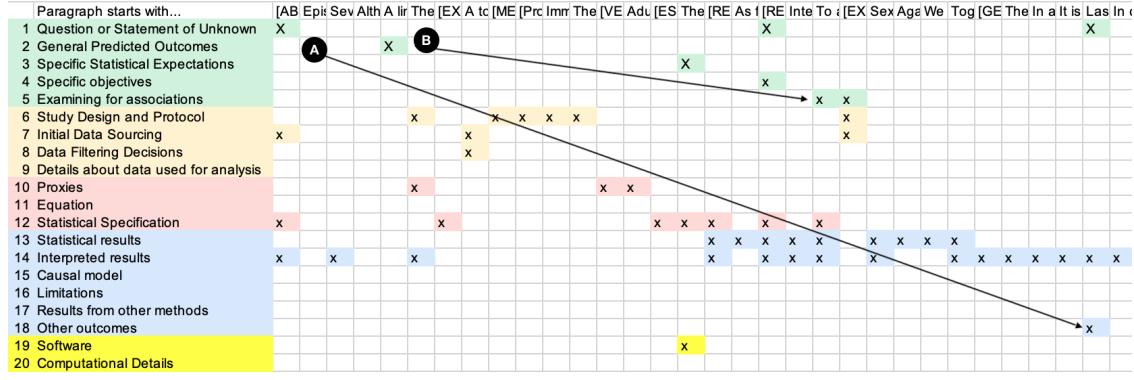


Figure 4.3: Formative content analysis: example reorderable matrix for [NHNO19].

We visualized each paper in our sample as a “reorderable matrix” [Ber11] to aid in detecting patterns in papers’ structure and content that could indicate how researchers formalized their hypotheses. The rows represent the codes in our codebook, colored according to the five broad categories of codes: research goals (rows 1-5, green), sample information (rows 6-9, orange), statistical analysis details (rows 10-12, red), reporting of results (rows 13-18, blue), and computational details (rows 19-20, bright yellow). The columns are the paragraphs, which are indexed by their first sentences, ordered left to right. In a paragraph’s column, there is an “X” for each code the paragraph received. Paragraphs have multiple codes if they contain multiple types of information. Among the ten visual patterns we noticed across our sample and subsequently looked for in each paper, two stand out in this paper. (A) As the paper progresses (visually moving left to right), the paper’s focus shifts from research goals to sample information to statistical analysis to results, as indicated by the arrow labeled A. Largely expected, this pattern helps to validate our coding method. Also, there is only one paragraph that discusses statistical software. (B) Researchers discuss research goals and questions throughout the paper. Interestingly, in the middle of the paper, when the researchers discuss their goals in greater detail, the researchers discuss them in increasing specificity, as indicated by the arrow labeled B. We were able to detect this pattern across papers by iterating on how to order the research goal codes (rows 1-5, green). The final order lists codes in increasing specificity from top (row 1) to bottom (row 5). Pattern B suggests that researchers refine their hypotheses during hypothesis formalization, which may involve specifying proxies and statistical methods. Chapter A discusses additional patterns in this paper and across our entire sample.

three authors then (i) read and coded all sections of papers except the figures, tables, and auxiliary materials that did not pertain to methodology²; (ii) discussed and summarized the papers’ goals and main findings to ensure comprehension and identify contribution types; and (iii) visualized each paper as a “reorderable matrix” [Ber11].

We adapted Bertin’s “reorderable matrix” [Ber11], an interactive visualization technique for data exploration, in our analysis. We visualized each paper in our sample as a matrix where each row represented a code in our codebook and each column represented a coded paragraph. We fixed the order of paragraphs to match the paper’s progression. We colored codes (rows) according to their categories in our codebook, repeatedly reordered the rows representing codes, and transposed the matrices to detect visual patterns in the papers. Figure 4.3 shows an example matrix.

²PNAS and Nature papers included a materials and methods section after references that were distinct from extended tables, figures, and other auxiliary material. We coded the materials and methods sections in the appendices and included them in the content analysis. Section A.2 describes our process in greater detail.

The visual representation of papers' content and structure helped us notice common patterns across papers and guided our follow-up analyses and discussions about what steps (**RQ1 - Steps**) and considerations (**RQ2 - Process**) researchers reported having during hypothesis formalization. Across multiple papers, the matrices showed how researchers typically start with broader research goals that they decompose into specific hypotheses (i.e., hypothesis refinement) over the course of a paper section, for example. Within a single paper, the matrices visually showed patterns of how researchers motivated and pieced together multiple experiments and interpreted statistical results in order to make a primary scientific argument. Chapter A includes our codebook with definitions and examples as well as a summary, citation, and annotated matrix for each paper.

Findings

Overview: We coded a total of 2,989 paragraphs across 50 papers. Results were the most commonly discussed topic. Approximately 31% of the paragraphs (in 50 papers) discussed interpretations of statistical results, and 11% (in 37 papers) provided details about statistical results (e.g., parameter estimates). Interpreted results often co-occurred with statistical results. 21% of paragraphs (in 40 papers) described data collection design (e.g., how the experiment was designed, how the data were collected, etc.). Specifications of statistical models appeared in 19% of paragraphs (in 50 papers). 11% of paragraphs (in 45 papers) discussed proxy variables, or measures to quantify abstract constructs (e.g., music enjoyment). To our surprise, more papers mentioned software than included equations. Researchers mentioned software used for statistical analysis in 3% of paragraphs (in 25 papers), sometimes even specifying function names and parameters, a level of detail we did not expect to find in publications. Only fifteen papers (JFE: 9, PS: 5, PNAS: 1) included equations in a total of 71 paragraphs. This suggests that mathematical equations, though part of the hypothesis formalization process, are less important to researchers than their tool-specific implementations.

We present more comprehensive tables and findings about paper structure, about paper contributions and venue differences in Chapter A.

Theme: Researchers decompose hypotheses into sub-goals that correspond to statistical analyses.

In approximately 70% of papers in the corpus, we found that researchers deconstructed their motivating research questions and overarching hypotheses into more tightly scoped objectives or relationships of interest that map to specific statistical analyses. For example, in [CDdME19], the researchers asked how theories of macroevolution varied across groups of species. The authors divided pre-existing hypotheses into three classes of hypotheses and assessed each class in turn. For one class of “geometric” hypotheses about insect egg size, the researchers discriminated between

two opposing hypotheses by examining “the scaling exponent of length and width (the slope of the regression of log-transformed length and log-transformed width).” As this example demonstrates, hypothesis formalization involves an *iterative hypothesis refinement process at the conceptual level*. This refinement process distills hierarchies of hypotheses and/or a single conceptual hypothesis into sub-hypotheses and formalizes these sub-hypotheses in statistical model implementations. Researchers also relate sub-hypotheses to one other during this process, which implies their causal models about the motivating conceptual hypothesis (and domain).

Theme: Researchers concretize hypotheses using proxies that are based on theory or available data.

Proxy variables further refine conceptual hypotheses by identifying how observable some concepts are, measuring the observable ones, indirectly measuring the less observable ones, and comparing measurement choices to other possible measures or ideal scenarios. As such, proxy variable selection is an important transition step between conceptual and data concerns during hypothesis formalization.

When defining proxy variables, researchers (i) used previously validated measures when available for theoretical and methodological soundness, such as the Barcelona Music Reward Questionnaire (BMRQ) to measure music reward (in [HAPGNN⁺19]), or (ii) developed new measures as a research contribution. For example, in [GAL⁺19], the authors develop an EEG-based measure for “immersiveness” in VR they demonstrated to be superior to previous measures that required halting immersive VR experiences to ask users about immersion. Researchers also sometimes justified choosing proxies based on available data. For example, in [GKM19], the researchers wanted to develop a proxy variable for job rank based on titles and “financial outcomes” (e.g., compensation, bonuses, etc.) to see if housing bankers were promoted or demoted after the 2008 stock market crash. However, because the financial outcomes were not public, the researchers relied on title only to compare bankers’ ranks, which was sub-optimal because job titles differ between companies.

Researchers consider their proxy choices as study limitations and consider alternative proxies to ensure that their findings are robust. Validating findings with multiple proxies suggests that hypothesis formalization can be a *recursive process*. Proxies lead to follow-up hypotheses about possible latent measurement factors, for instance, which in turn lead to additional analyses that address the same conceptual hypothesis.

Theme: Data collection and sampling influence statistical analysis.

Researchers often described their data sampling and study design as factors that necessitated additional steps in their analysis process. In [BLT⁺19] and [PAH⁺19], researchers accounted for effects

of task order in their study protocol by running additional regressions or analyzing tasks separately. Researchers also ran initial analyses to assess the effect of possibly confounding variables in their study design, such as gender in [CLD⁺19] or location of stimuli in [TPO⁺19]. Other times, researchers performed robustness checks after their main analyses, such as in response to a gender imbalance in [PAH⁺19] and possible sample selection biases due to database constraints in [HHZ19].

Although data collection driven by statistical modeling plans was expected of replication studies (e.g., [ZTDB19; PAH⁺19; BLT⁺19]) or papers that make methodological contributions (e.g., [CCM19a; CCM19b]), we found an instance in [BSCN19]—neither replication nor methodological contribution—where researchers explicitly reported selecting a statistical model before designing their study. The researchers chose to use a well-validated computational model, the linear ballistic accumulator (LBA), to quantify aspects of human decision making. This model selection influenced the way they designed their study protocol so that they could obtain a sample large enough for accurate parameter estimation.

Based on these observations, it seems that modeling choices more frequently react to data collection processes and possible sample biases, following a linear data collection-first process implied by prior work. However, there are also instances where model implementation comes first and researchers' data collection procedures must adhere to modeling needs, suggesting a previously missing *loop between statistical model implementation and data collection* that is separate from any influences conceptual hypotheses have on data collection.

4.2.1 Discussion

The content analysis confirmed prior findings on (i) the connection between hypotheses and causal models (e.g.,[McE20]), (ii) the importance of proxies to quantify concepts, and (iii) the constraints that data collection design and logistics place on modeling. Extending prior work, the content analysis also (i) suggested that decomposing hypotheses into specific objectives is a mechanism by which conceptual hypotheses relate to causal models; (ii) crystallized the hypothesis refinement loop involving conceptual hypotheses, causal models and proxies; and (iii) surfaced the dual-search nature of hypothesis formalization by suggesting that model implementation may shape data collection.

The content analysis also raised questions about how much the pressures to write compelling scientific narratives [Ker98] influence which aspects of hypothesis formalization are described or omitted (e.g., in practice, model implementations may constrain data collection more often than we found in our dataset), how the steps are portrayed linearly even though the process may have been more iterative, how analysts determine which tools to use, and how analysts without domain expertise may approach hypothesis formalization differently. These questions motivated us to design and conduct a lab study to provide greater visibility into how analysts who are not necessarily

researchers approach the process with expectations of rigor but without pressure of publication.

Limitations

The major limitation of analyzing published papers is the disconnect between actual and reported analytical practice. The pressures to write compelling scientific narratives [Ker98] likely influence which aspects of hypothesis formalization are described or omitted. For instance, in practice, model implementations may constrain data collection more often than we found in our sample. Nevertheless, the lack of information in prior work and the content analysis suggests that hypothesis formalization remains an opaque process deserving of greater scrutiny. Hypothesis formalization may explain how analysts determine which tools to use and how domain expertise may influence the analytical conclusions reached.

4.2.2 Takeaways: Expected Steps in Hypothesis Formalization

Towards our first two research questions about what actions analysts take to formalize hypotheses (**RQ1 - Steps**) and why (**RQ2 - Process**), prior work and our formative content analysis suggest that hypothesis formalization involves steps in three categories: conceptual, data-based, and statistical. *Conceptually*, analysts develop conceptual hypotheses and causal models about their domain that guide their data analysis. With respect to *data*, analysts explore data and incorporate insights from exploration, which can be top-down or bottom-up, into their process of formalizing hypotheses. The *statistical* concerns analysts must address involve mathematical and computational concerns, such as identifying a statistical approach (e.g., linear modeling), representing the problem mathematically (e.g., writing out a linear model equation), and then implementing those using software. In our work, we find evidence to support separating statistical considerations into concerns about mathematics, statistical specification in tools, and model implementation using tools.

A key observation about prior work is that there is a tension between iterative and linear workflows during hypothesis formalization. Although sensemaking processes involve iteration, concerns about methodological soundness, as evidenced in pre-registration efforts that require researchers to specify and follow their steps without deviation, advocate for, or even impose, more linear processes. More specifically, theories of sensemaking that draw on cognitive science, in particular [RSPC93; GW14], propose larger iteration loops between conceptual and statistical considerations. Some textbooks and research concerning statistical thinking and practices [WP99; CEG⁺16] appear less committed to iteration while other researchers and practitioners in applied statistics emphasize *workflows* for iterating on statistical models [YK20; LCD⁺19a; GCS⁺13]. Workflows (e.g., model expansion) can help researchers start with simple models and build up to more complex ones by incrementally testing and refining their understanding of characteristics of the data,

the model fitting algorithms, and computational settings [Bet20; GVS⁺20a; GSV⁺19]. Moreover, empirical work in HCI on data analysis embraces iteration during exploration and observes iteration during some phases of confirmatory data analysis, such as statistical model choice, but not in others, such as tool selection. In our work, we are sensitive to this tension and aim to provide more granular insight into iterations and linear processes involved in hypothesis formalization. We also anticipate that the steps identified in prior work will recur in our lab study, but we do not limit our investigation to these steps.

4.3 Exploratory Lab Study

To address the limitation of the content analysis, understand analysts' considerations (**RQ2 - Process**) while formalizing their hypotheses (**RQ1 - Steps**), and examine the role of statistical software in this process (**RQ3 - Tools**), we designed and conducted a virtual lab study with freelance data workers who approach the hypothesis formalization and analysis process with expectations of rigor but without the pressure of publication.

4.3.1 Methods

Data workers: We recruited 24 data workers with experience in domains ranging from marketing to physics to education through Upwork (22) and by word of mouth (2).³

Twelve data workers held occupations as scientists, freelance data scientists, project managers, or software engineers. Six were currently enrolled in or had just finished graduate programs that involved data analysis. Five identified as current or recent undergraduates looking for jobs in data science. One was an educator. Data workers self-reported having significant experience on a 10-point scale adapted from a scale for programming experience [FKL⁺12] (min=2, max=10, mean=6.4, std=2.04) and would presumably have familiarity with hypothesis formalization.

The lab study enables us to contrast normative expert practices (found in prior work and our formative content analysis) to observed practices with data workers who are not statistical experts but still work in real-world analysis settings (i.e., research, marketing, consulting). A benefit of studying these data workers is that they are likely to benefit most from new tools.

Protocol: We designed and conducted a lab study with three parts. Parts 1 and 3 were recorded and automatically transcribed using Zoom. We compensated data workers \$45 for their time. The first author conducted the study and took notes throughout.

Part 1: Structured Tasks. Part 1 asked data workers to imagine they were leading a research

³We refer to our participants as data workers because they work with data but do not represent the entire population of data scientists, which may include statistical experts.

team to answer the following research question: “What aspects of an individual’s background and demographics are associated with income after they have graduated from high school?”⁴ We asked data workers to complete the following tasks:

- *Task 1: Hypothesis generation.* Imagining they had access to any kind of data thinkable, data workers brainstormed at least three hypotheses related to the research question.
- *Task 2: Conceptual modeling.* Next, data workers saw a sample data schema and developed a conceptual model for one or more of their hypotheses. We used the term “conceptual model” instead of “causal model” to avoid (mis)leading data workers. We provided the following definition: “A conceptual model summarizes the process by which some outcome occurs. A conceptual model specifies the factors you think influence an outcome, what factors you think do not influence an outcome, and how those factors might interact to give rise to the outcome.”
- *Task 3: Statistical model specification.* Finally, we presented data workers with a sample dataset and instructed them to specify but not implement a statistical model to test one or more of their hypotheses.

After the three tasks, we conducted a semi-structured interview with data workers about (i) their validity concerns⁵ and (ii) experiences. To help us contextualize our observations and assess the generalizability of our findings, we asked data workers to compare the study’s structure and tasks to their day-to-day data analysis practices.

Part 2: Take-home analysis. After the first Zoom session, data workers implemented their analyses using the previously shown dataset, shared any analysis artifacts (e.g., scripts, output, visualizations, etc.), and completed a survey about their implementation experience. Prior to Part 3, the first author reviewed all submitted materials and developed participant-specific questions for the final interview.

Part 3: Final Interview. The first author asked data workers to give an overview of their analysis process and describe the hypotheses they tested, how their analysis impacted their conceptual model and understanding, why they made certain implementation choices, what challenges they faced (if any), and any additional concerns about validity.

Materials: The data schema and dataset used in the study came from a publicly available dataset from the Pew Research Center [Suh14]. Each task was presented in a separate document. All study materials are included in the appendix.

⁴We chose the open-ended research question about income after high school because we expected it to be widely approachable and require no domain expertise to understand.

⁵If data workers were unfamiliar with the term “validity,” we rephrased the questions to be about “soundness” or “reliability.”

Analysis: The first author reviewed the data workers' artifacts multiple times to analyze their content and structure; thematically analyzed notes and transcripts from data workers' Zoom sessions; and regularly discussed observations with the other authors throughout analysis.

4.3.2 Findings and Discussion

Eighteen of the 24 data workers we recruited completed all three parts of the study. The other six data workers completed only the first Zoom session. In our analysis, we incorporate data from all data workers for as far as they completed the study.

We found that data workers had four major steps (**RQ1 - Steps**) and considerations (**RQ2 - Process**): (i) identifying or creating proxies, (ii) fitting their present analysis to familiar approaches, (iii) using their tools to specify models (**RQ3 - Tools**), and (iv) minimizing bias by relying on data. Data workers also faced challenges acquiring and incorporating domain and statistical knowledge (**RQ2 - Process**).

Theme: Data workers consider proxies and data collection while articulating hypotheses.

We encouraged data workers to not consider the feasibility of collecting data while brainstorming hypotheses. Yet, while brainstorming hypotheses, data workers expressed concern with how to measure constructs [D2, D5, D8, D12, D18, D22, D24] and how to obtain data [D2, D6, D8, D9, D11, D21, D24].

For instance, D18, a computer science student who had worked on more than five data analysis projects, grappled with the idea of 'privilege' and how to best quantify it:

"I'm trying to highlight the fact that those who will be privileged before graduation...that experience will enable them to make again more money after graduation. I won't say 'privilege' because we need to quantify and qualify for that...it's just an abstract term."

Eventually, D18 wrote two separate hypotheses about 'privilege,' operationalizing it as parental income: (1) "People with higher incomes pre graduating, end up having higher differences between pre and post graduation incomes than those with lower incomes pre graduation." and (2) "People with parents with lower incomes tend to have lower incomes pre graduation than those with parents with higher incomes."

D18 continued to deliberate 'privilege' as measured by low and high income, saying, "...again you need to be careful with low and high because these are just abstract terms. We need to quantify that. What does it mean to be 'low?' What does it mean to be 'high?'". Finally, D18 decided to "maybe use the American standards for low income and high income." Although an accepted "American

standard” may not exist, D18 nevertheless believed that cultural context was necessary to specify because it could provide a normalizing scale to compare income during analysis, demonstrating how data workers plan ahead for statistical modeling while brainstorming and refining hypotheses.

Similarly, D2, a freelance data scientist, was very specific about how to measure personality: “More extraverted individuals (extraversion measured using the corresponding social network graph) are likely to achieve higher yearly income later in life.”

In the presence of the data schema, more data workers were concerned with proxies [D2, D5, D6, D7, D8, D9, D16, D18, D21]. Some even adapted their working definitions to match the available data, similar to how researchers in the content analysis determined proxies based on data. For instance, D8, who hypothesized that “individuals interested in STEM fields tend to earn more post high school than individuals interested in other fields,” operationalized “interest” as “Major” — a variable included in the data schema — even though they had previously brainstormed using other proxies such as club attendance in high school.

These data workers’ closely related considerations of data and concept measurement demonstrate how conceptual hypotheses and data collection may inform each other, corroborating our findings from the content analysis.

Theme: Data workers consider implementation and tools when specifying statistical models.

When we asked data workers to specify their models without considering implementation, we anticipated they would name specific statistical tests (e.g., “ANOVA”), approaches (e.g., “linear regression” or “decision trees”), or write mathematical models (e.g., $Y = B_0 + B_1X_{age} + B_2X_{gender}$) that they could then implement using their tools because (a) some researchers in the literature survey did so in their papers and (b) several data workers mentioned having years of analysis experience. However, despite the explicit instruction to disregard implementation, 16 data workers provided to-do lists or summaries of steps to perform a statistical analysis as their model specifications [D1, D2, D3, D5, D7, D8, D9, D11, D12, D14, D16, D18, D20, D21, D22, D23, D24]. Of these 16 data workers, eight also named specific statistical tests in their descriptions [D3, D7, D8, D11, D12, D14, D18, D20].

For example, D8, a data science consultant with 7/10 analysis experience, specified a list of steps that included creating new variables that aggregated columns in the dataset, cleaning and wrangling the data, visualizing histograms, performing chi-squared test, and interpreting the statistical results. Notably, D8 also specified null and alternative hypotheses, which acted as an intermediate artifact during hypothesis formalization. Figure 4.4 shows D8’s statistical specification.

Only four data workers named specific statistical methods without describing their steps [D4, D6, D15, D17]. Two data workers, D22, a neuroscientist by training with 8/10 analysis experience, and

Create new variables:

Adj_annual_income - take the midpoint of the ranges in the Annual Income column as a numeric value. (numeric)

State_avg_income - find the average income of individuals in each state from established benchmarks. (numeric)

Income_over_avg - take the difference between each individual's income with the average for their state.

Testing Major vs income: take all rows with a college degree (2 year associate and up) & major. Omit rows with no info on income.

For each major, calculate the average *Adj_annual_income*.

Also, calculate the average *Adj_annual_income* for all the college rows from above.

Create a set of histograms (one for each major) showing the spread of *Adj_annual_income* for the people in that group. The histograms should share the same x axis. The bins will be normalized to sum to 100% for each major group.



Arrange the data like so

| Major | Avg Income (within major) | Avg income (sample population) |
|-------|---------------------------|--------------------------------|
| Bio | #### | #### |
| Stats | #### | #### |
| etc. | #### | #### |

Chi-squared test.

H₀: for each major group, the average income is equal to the entire sample population's average income. That is, no single group has a significant difference in avg income from the sample population.

H_A: at least one of the major groups has an average income that's significantly different from the sample population.

Test for a p-value <= 0.05

One caveat of our selected test is even if we are able to reject H₀, we can't make conclusions about which major group is the one making the different. It's possible that just one group is; it's possible that every group is significantly different from the population wrt large.

Figure 4.4: Sample statistical specification (D8).

The lab study tasked analysts to specify their statistical models without considering implementation. We expected analysts would represent their statistical models using statistical test names or mathematical equations. Instead, most analysts specified statistical procedures for performing statistical models using todo lists and summaries of steps, which sometimes included mentions of software tools, showing that implementation was an important consideration and that tool familiarity may limit which statistical models analysts consider and implement. Data worker D8 specified their model through a combination of statistical test names (e.g., Chi-squared test) and a list (split across two pages) of detailed steps involved in creating new variables, cleaning and wrangling data, visualizing data, and testing their hypothesis.

D19, an educator with 6/10 analysis experience, attempted to specify their models mathematically. D22 used the familiar R syntax: “Current Income ~ Educational attainment + Gender + Interactions of those two.” On the other hand, D19 gave up because although they knew the general form of logistic regression, they did not know how to represent the specific variables in the model they wanted to perform.

The implementation and software details data workers discussed and included in their specifications suggest that data workers prefer to skip over mathematical equations and jump to specification and implementation in their tools. Although it is possible that study instructions primed data workers to respond about how they would perform, rather than represent, the task even after researcher clarifications, this would not explain the level of implementation detail data workers included. Nine data workers went so far as to mention specific libraries, even functions, that they would use to program their analyses [D3, D9, D12, D13, D14, D16, D19, D21, D23]. In their reflective interviews, data workers also expressed that they often do not specify models outside of implementing them, which D19 succinctly described:

“I don’t normally write this down because all of this is in a [software] library.”

Data workers’ statistical knowledge appears to be situated in the programs they write, and their knowledge of and familiarity with tools constrains the statistical methods they explore and consider. As such, tools may be a key point of intervention for guiding data workers toward statistical methods that may be unfamiliar but are best suited for their conceptual hypotheses.

Theme: Data workers try to fit analyses to previous projects and familiar approaches.

Data workers spent significant thought and time categorizing their analyses as “prediction,” “classification,” or “correlation” problems [D2, D3, D7, D10, D11, D18, D19, D21, D22]. To categorize, data workers relied on their previous projects. While reflecting on their typical analysis process, D21, a software engineer working in healthcare, said (emphasis added),

*“I usually tend to jump...to look at data and **match** [the analysis problem] with similar patterns I have seen in the past and start implementing that or do some rough diagrams [for thinking about parameters, data type, and implementation] on paper...and start implementing it.”*

Data workers also looked at variable data types (i.e., categorical or continuous) to categorize. For example, D3, a freelance analyst, pivoted from thinking about **predicting** income to **classifying** income groups (emphasis added) based on data type information:

*“The income, the column, the target value here, is categorical. I think maybe it wouldn’t be a bad idea to see what **classification** tasks, what we could do. So instead of trying*

to predict because we're not trying to predict an exact number, it seems...like more of a classification problem..."

A provocative case of adhering to prior experiences was D6, a psychological research scientist. Although several data workers were surprised and frustrated that income was ordinal in the dataset with categories such as "Under \$10K," "\$10K to \$20K," "\$20K to \$30K," up to "150K+", none went so far as D6 to synthetically generate normally distributed income data so that they could implement the linear regression models they had specified despite saying they knew that income was not normally distributed.

When asked further about the importance of normal data, D6 described how they plan analyses based on having normal data, strive to collect normally distributed, and rely on domain knowledge to transform the data to be normal when it may not be after collection:

"...I feel like having non normal data is something that's like hard for us to deal with. Like it just kind of messes everything up like. And I know, I know it's not always assumption of all the tasks, but just that we tend to try really hard to get our variables to be normally distributed. So, you know, we might like transform it or, you know, kind of clean it like clean outliers, maybe transform if needed...I mean, it makes sense because like a lot of measures we do use are like depressive symptoms or anxiety symptoms and kind of they're naturally normally distributed...I can probably count on my hand the number of non parametric tests I've like included in manuscripts."

D6's description of their day-to-day analyses exemplifies the dual-search nature of hypothesis formalization: Data workers (i) jump from hypothesis refinement to model specification or implementation with specific proxies in mind and then (ii) collect and manipulate their data to fit their model choices.

We recognize that data workers may have taken shortcuts for the study they would not typically make in real life. Nevertheless, the constraints we imposed by using a real-world dataset are to be expected in real-world analyses. Therefore, our observations still suggest that rather than consider the nature and structure of their hypotheses and data to inform using new statistical approaches, which statistical pedagogy and theory may suggest, data workers may choose familiar statistical approaches and mold their new analyses after previous ones.

Theme: Data workers try to minimize their biases by focusing on data.

Throughout the study, data workers expressed concern that they were biasing the analysis process. Data workers drew upon their personal experiences to develop hypotheses [D5, D10, D13, D15, D16, D20, D21, D24] and conceptual models [D8, D12, D20, D24]. D12, a data analysis project manager,

described how their personal experiences may subconsciously bias their investigation by comparing a hypothetical physicist and social worker answering the same research question:

“Whereas a social worker by design...they’re meant to look at the humanity behind the numbers [unlike a physicist]. So like, they may actually end up with different results...actually sitting in front of this data, trying to model it.”

A few data workers even refused to specify conceptual models for fear of biasing the statistical analyses [D10, D11, D19]. On the surface, data workers resisted because they believed that some relationships, such as the effect of age on income, were too “obvious” and did not warrant documentation [D10, D11]. However, relationships between variables that were “obvious” to some data workers were not to others. For instance, D10, a business analyst, described how income would plateau with age, but other data workers, such as D18, assumed income would monotonically increase with age.

When we probed further into why D10, D11, and D19 rejected a priori conceptual models, they echoed D10’s belief that conceptual models “put blinders on you.” Even the data workers who created conceptual models echoed similar concerns of wanting to “[l]et the model do the talking” in their implementations [D3, D15, D18, D19]. Instead of conceptual modeling, D10 chose to look at all n-ary relationships in the dataset to determine which variables to keep in a final statistical model, saying,

“It’s so easy to run individual tests... You can run hypothesis tests faster than you can actually think of what the hypothesis might be so there’s no need to really presuppose what relationships might exist [in a conceptual model].”

Of course, one could start from the same premise that statistical tests are so easy to execute and conclude that conceptual modeling is all the more important to prioritize analyses and prevent false discoveries.

Similarly, data workers were split on whether they focused their implementation exclusively on their hypotheses or examined other relationships in the dataset opportunistically. Nine data workers stuck strictly to testing their hypotheses [D1, D4, D5, D6, D7, D11, D13, D20, D24]. However, five data workers were more focused on exploring relationships in the dataset and pushed their hypotheses aside [D2, D3, D10, D16, D18], and an additional four data workers explored relationships among variables not previously specified in their hypotheses in addition to their hypotheses [D14, D15, D17, D21]. D18 justified their choice to ignore their hypotheses and focus on emergent relationships in the data by saying that they wanted to be “*open minded based on the data...open to possibilities.*”

Data workers' concerns about bias and choice of which relationships to analyze (hypothesis only vs. opportunistic) highlight the tension between the two searches involved in hypothesis formalization: concept-first model implementations and implementation-first conceptual understanding. Conceptual models are intermediate artifacts that could reconcile the two search processes and challenge data workers' ideas of what "data-driven" means. However, given some data workers' resistance to prior conceptual modeling, workflows that help data workers conceptually model as a way to reflect on their model implementations and personal biases may be more promising than ones that require them before implementation.

Theme: Data workers face challenges obtaining and integrating conceptual and statistical information.

Based on data workers' information search behaviors and self-reports, we found that data workers faced challenges obtaining and integrating both domain and statistical knowledge.

Data workers consulted outside resources such as API documentation, Wikipedia, and the *Towards Data Science* blog throughout the study: one while brainstorming hypotheses [D13]; three while conceptual modeling [D12, D13, D22]; six while specifying statistical models [D3, D6, D12, D13]. Six data workers also mentioned consulting outside resources while implementing their analyses [D1, D3, D11, D14, D15, D21]. By far, statistical help was the most common.

Furthermore, when data workers reflected on their prior data analysis experiences, they detailed how collaborators provided domain and statistical expertise that are instrumental in formalizing hypotheses. Collaborators share data that help domain experts generate hypotheses [D9], critique and revise conceptual models and proxies [D4, D8], answer critical data quality questions [D10], and ensure statistical methods are appropriate [D5, D6, D22].

In the survey participants completed after implementing their analyses, the three most commonly reported challenges were (i) **formatting** the data [D1, D4, D5, D6, D13, D16, D18, D20, D21, D24], (ii) **identifying** which statistical analyses to perform with the data to test their hypotheses [D1, D11, D14, D18, D20, D21], and (iii) **implementing and executing** analyses using their tools [D1, D6, D7, D13, D20, D21]. Although we expected data workers would have difficulty wrangling their data based on prior work [KPHH12], we were surprised that identifying and executing statistical tests were also prevalent problems given that (a) data workers were relatively experienced and (b) could choose their tools. These results, together with our observations that data workers rely on their prior experiences and tools, suggest that data workers have difficulty adapting to new scenarios where new tools and statistical approaches may be necessary.

4.3.3 Takeaways from the Lab Study

After the first session, 13 out of the 24 data workers described all the tasks as familiar, and 10 described most of the tasks and process as familiar. Data workers commonly remarked that although the process was familiar, the order of the tasks was “opposite” of their usual workflows. In practice, data workers may start with model implementation before articulating conceptual hypotheses, which opposes the direction of data analysis that the ASA recommends [CEG⁺16]. Nevertheless, our observations reinforce the dual-search, non-linear nature of hypothesis formalization.

Moreover, one data worker, D24, a physics researcher who primarily conducted simulation-based studies expressed that the study and its structure felt foreign, especially because they had no control over data collection. Other data workers in the study also described the importance of designing and conducting data collection as part of their hypothesis formalization process [D4, D6, D9]. Designing data collection methods informs the statistical models data workers plan to use and helps to refine their conceptual hypotheses by requiring data workers to identify proxies and the feasibility of collecting the proxy measures, reinforcing what we saw in the content analysis. The remarks also suggest that disciplines practice variations of the hypothesis formalization process we identify based on discipline-specific data collection norms and constraints. For example, simulating data may sometimes take less time than collecting human subjects data, so data workers working with simulations may dive into modeling and data whereas others may need to plan experiments for a longer period of time.

Approximately half of the data workers had either just finished or were enrolled in undergraduate or graduate programs involving data analysis. As such, half of our sample likely has limited professional experience outside of their studies and/or freelance work on Upwork. Additionally, data work available on Upwork may be more narrowly focused and less representative of end-to-end data analysis or research projects expected of those with greater statistical expertise. Still, several data workers in our study mentioned other employments where they gained professional experience working on larger analysis and research projects. Despite the limitations of recruiting participants from Upwork and word of mouth, our sample represents data workers who have training in a diversity of disciplines (e.g., medicine, psychology, business), are familiar with a range of statistical methods, and have experience using a broad range of statistical tools. As such, the data workers in our study may be representative of analysts who are likely to benefit most from new tools for supporting hypothesis formalization.

Finally, we found that data workers relied on prior experiences and tools to specify and formalize their hypotheses. Tools that scaffold the hypothesis formalization process by suggesting statistical models that operationalize the conceptual hypotheses, conceptual models, or partial specifications data workers create along the way may (i) nudge data workers towards more robust analyses that

test their hypotheses, (ii) overcome limitations of data workers’ prior experiences, and (iii) even expand data workers’ statistical knowledge. Thus, we investigated how current tool designs serve (or under-serve) hypothesis formalization.

4.4 Analysis of Software Tools

To understand how the design of statistical computing tools may support or hinder hypothesis formalization (**RQ3 - Tools**), we analyzed widely used software packages and suites. Throughout, we use the term “package” to refer to a set of programs that must be invoked through code, such as `lme4`, `scipy`, and `statsmodels`. We use the term “suite” to refer to a collection of packages that end-users can access either through code or graphical user interfaces (GUIs), such as SPSS, SAS, and JMP. We use the term “tool” to refer to both. Software packages were a unit of analysis because they are necessary for model implementation regardless of medium (e.g., computational notebook, CoLab, RStudio). As such, our findings apply to tools that provide wrappers around packages included in our sample.

4.4.1 Method

Sample: Our sampling procedure involved two phases: (i) identifying software packages and suites for model implementation (not visual analysis tools like Tableau) mentioned more than once across the content analysis and lab study and (ii) adding recommended packages and suites from online data science communities our lab participants mentioned or used (e.g., *Towards Data Science*). To identify these additional tools, we consulted online data analysis fora [Gro19; Bob17; Bob18; Pra19]. The final sample included 20 statistical tools: 14 packages (R: 10, Python: 4); three suites that support in-tool programming; and three suites that do not support programming. Table 4.1 contains an overview of our sample and results.

Analysis: Four specific questions guided our analysis:

- **Specialization:** Data workers in the lab study eagerly named specific statistical tools they would use and looked up tool documentation during the tasks. This prompted us to ask, *How specialized are the tools, and how might specialization (or lack thereof) affect how end-users discover and use them to formalize hypotheses?*
- **Statistical Taxonomies:** Data workers in the lab study tried to mold their analyses to prior experiences and their taxonomies of statistical methods. We wondered what role tools play in this: *How do tools organize and group statistical models? How might tool organization and end-users’ taxonomies interplay during hypothesis formalization?*

- **Model Expression:** Data workers in the lab study jumped to model implementation throughout the tasks. Only half provided names of statistical methods. We wondered if this was due to how tools enable end-users to express their models: *What notation must end-users use to express models in the tools?*
- **Computational Issues:** Data workers in the lab study described their statistical models using specific function calls. Similarly, although it was uncommon for researchers in the content analysis to specify the software tools they used, when they did, researchers specified the functions, parameters, and settings used. This prompted us to wonder about the importance of computational settings: *What specific kinds of computational control do tools provide end-users and how might that impact hypothesis formalization?*

To answer the four questions for each statistical tool, the first author read and took notes on published articles about tools' designs and implementations, API documentation and reference manuals, and available source code; followed online tutorials; consulted question-and-answer sites (e.g., StackExchange) when necessary; and analyzed sample data with the tools. The first author paid particular attention to tool organization, programming idioms, functions and their parameters, and tool failure cases. Table 4.1 contains citations for resources consulted in the analysis. The iterative analysis process involved discussions among the co-authors about how to evaluate the properties of tools from our perspectives as both tool designers/maintainers and end-users. Here, we focus on end-user (hereafter referred to as analyst) perspectives informed by our lab study and make callouts to details relevant for tool designers.

4.4.2 Findings and Discussion

We discuss our findings in light of our characterization of hypothesis formalization in Figure 4.1. We refer to specific steps and transitions in Figure 4.1 in **boldface**.

Theme: Specialization.

Half the tools [T2, T3, T4, T5, T6, T7, T8, T9, T11, T12] in our sample are specialized in the scope of statistical analysis methods they support (e.g., `brms` supports Bayesian generalized linear multilevel modeling). `edgeR` [T3] provides multiple modeling methods but is specialized to the context of biological count data. Such specialized tools are vital to creating a widely adopted statistical computing ecosystem, such as R.

Despite its importance, tool specialization pushes computational concerns higher up the hypothesis formalization process. Specialized tools require analysts to consider computational settings while picking a statistical tool and, possibly, even while mathematically relating their variables. They fuse

Table 4.1: Overview of the software tools included in our analysis.

Half of the tools are specialized for specific modeling use cases. Most tools use mathematical notation (T18–T20 (\checkmark^*) even use mathematical notation in their GUIs). Most tools also provide a wide range of computational control although sometimes they require additional packages [T5, T13]. Tool specialization, organization, notation, and computational control focus analysts on model implementation details, sometimes at the expense of focusing on their conceptual hypotheses.

| ID | Tool name | Specialized Scope | Mathematical Notation | Computational Control | References |
|--|------------------------------------|-------------------|-----------------------|-----------------------|--|
| R Packages | | | | | |
| T1 | MASS | — | ✓ | ✓ | [RVB ⁺ 20] |
| T2 | brms | ✓ | ✓ | ✓ | [B ⁺ 17; BB16] |
| T3 | edgeR | ✓ | ✓ | ✓ | [CLM ⁺ 20; CMR ⁺ 20] |
| T4 | glmmTMB | ✓ | ✓ | ✓ | [BKvB ⁺ 17; MSN ⁺ 20] |
| T5 | glmnet | ✓ | — | ✓(additional) | [FHT ⁺ 20; HQ14] |
| T6 | lme4 | ✓ | ✓ | ✓ | [BMBW14; BMB ⁺ 19] |
| T7 | MCMCglmm | ✓ | ✓ | ✓ | [H ⁺ 10; Had20] |
| T8 | nlme | ✓ | ✓ | ✓ | [PBD ⁺ 20] |
| T9 | RandomForest | ✓ | ✓ | ✓(minimal) | [BCLW18] |
| T10 | stats (core library) | — | ✓ | ✓ | [Tcw20] |
| Python Packages | | | | | |
| T11 | Keras | ✓ | — | ✓(minimal) | [C ⁺ 15] |
| T12 | Scikit-learn | ✓ | — | ✓ | [sld20; PVG ⁺ 11; BLB ⁺ 13] |
| T13 | Scipy (scipy.stats) | — | — | ✓(additional) | [JOP ⁺ 21a; JOP ⁺ 21b; JOP ⁺ 21c] |
| T14 | Statsmodels | — | ✓ | — | [SP10; PSTsd20] |
| Suites, with DSLs for programming | | | | | |
| T15 | Matlab (Statistics and ML Toolbox) | — | — | ✓ | [TM20a; TM20b] |
| T16 | SPSS | — | ✓ | ✓ | [SPS21] |
| T17 | Stata | — | ✓ | — | [Sta21; LLC20b; LLC20a] |
| Suites, without programming | | | | | |
| T18 | GraphPrism | — | ✓* | ✓ | [GS20] |
| T19 | JASP | — | ✓* | — | [oA20] |
| T20 | JMP | — | ✓* | — | [SAS20a; JS11] |

the last two steps of hypothesis formalization (**Statistical Specification** and **Model Implementation**). Ultimately, specialization requires analysts to have more (i) computational knowledge and (ii) foresight about their model implementations at the cost of focusing on conceptual or data-related concerns early in hypothesis formalization.

One way tool designers minimize the requisite computational knowledge and foresight while providing the benefits of specialized packages — which may be optimal for specific statistical models or data analysis tasks — is to provide micro-ecosystems of packages. For example, R’s `tidymodels` [KW20] and `tidyverse` [WAB⁺19] create micro-ecosystems that use consistent API syntax and semantics across interoperable packages. They also push analysts towards what the tool designers believe to be best practices, such as the use of the tidy data format [W⁺14]. Tools that aim to support hypothesis formalization may consider fitting into or creating micro-ecosystems that provide tool support all along the process, focusing analysts on concepts, data, or model implementation at various points.

Theme: Statistical taxonomies.

A consequence of tool specialization is the fragmented view of statistical approaches. For example, we observed analysts in the lab study who viewed the analysis as a classification task gravitate towards machine learning-focused libraries, such as `RandomForest` [T9], `Keras` [T11], and `scikit-learn` [T12]. Because classification can be implemented as logistic regression, any tool that supports logistic regression, such as the core `stats` library in R [T10], provides equally valid, alternative perspectives on the same analysis and hypothesis. However, tools obfuscate these connections and do not aid analysts in considering reasonable statistical models that may be unfamiliar or outside their personal taxonomy. This may explain why analysts adhered to their personal taxonomies during the lab study.

This problem carries over to tools that support numerous statistical methods. Ten tools in our sample intend to provide more comprehensive statistical support [T1, T10, T13, T14, T15, T16, T17, T18, T19, T20]. These tools group statistical approaches using brittle and inconsistent taxonomies based on data types [T17]; analysis classes that are both highly specific (e.g., “Item Response Theory”) and vague (e.g., “Multivariate analyses”) [T15, T16, T17, T18, T19, T20]; and disciplines or applications (e.g., “Epidemiology and related,” “Direct Marketing”) [T16, T17, T20]. Although well-intended to simplify statistical method selection, tools’ taxonomies are at times misleading. For instance, JMP combines various linear models into a “Fit Model” option that is separate from “Predictive Modeling” and “Specialized Modeling,” which are also distinct from the more general “Multivariate Methods.” Once analysts select the “Fit Model” option, they can specify the “Personality” of their model as “Generalized Regression,” “Generalized Linear Model,”

or “Partial Least Squares,” among many others. This JMP menu structure implies that (i) a Partial Least Squares model is distinct from a regression model when it is in fact a type of regression model and (ii) regression is not useful for prediction, which is not the case.

In these ways, tools add a “Navigate taxonomies” step before the **Statistical Specification** step, requiring analysts to match their conceptual hypotheses with the tools’ taxonomies, which may misalign with their personal taxonomies. One reason for this issue may be that tools do not leverage analysts’ intermediate artifacts or understanding during hypothesis formalization. By the time analysts transition to **Statistical Specification**, they have refined their conceptual hypotheses, developed causal models, and made observations about data. However, tools’ taxonomies require analysts to set these aside and consider another set of decisions imposed by tool-specific groupings of statistical methods. In this way, tool taxonomies may introduce challenges that detract from hypothesis formalization.

Theme: Model expression: Syntax and semantics

Fifteen tools in our sample provide analysts with interfaces that use mathematical notation to express statistical models [T1, T2, T3, T4, T6, T7, T8, T9, T10, T14, T16, T17, T18, T19, T20]. R and Python packages use symbolic mathematical syntax, and SPSS and Stata use natural language-like syntax. Expressing a linear model with Sex, Race, and their interaction as predictors of Annual Income involves the formula `AnnualIncome ~ Sex + Race + Sex*Race` in `lme4` and `AnnualIncome BY Sex Race Sex*Race` in SPSS. In a linear execution of steps involved in hypothesis formalization where analysts relate variables mathematically (**Mathematical Equation**) before specifying and implementing models using tools (**Statistical Specification**, **Model Implementation**), the mathematical interfaces match analysts’ progression. However, in the lab study, analysts did not specify their models mathematically even when given the opportunity, suggesting that mathematical syntax may not adequately capture analysts’ conceptual or statistical considerations.

Syntactic similarity between packages may lower the barrier to trying and adopting new statistical approaches that more directly test hypotheses and therefore benefit hypothesis formalization. At the same time, syntactic similarity may also introduce unmet expectations of semantic similarity. For example, `brms` [T2] uses the same formula syntax as `lme4` [T6], smoothing the transition between linear modeling and Bayesian linear modeling for analysts. However, based on syntactic similarity, analysts may incorrectly assume statistical equivalence in computed model values. For example, in `brms`, the model intercept is the mean of the posterior when all the independent variables are at *their means*, but in `lme4`, the intercept is the mean of the model when all the independent variables are at *zero*.

Conversely, tools introduce syntactic differences between statistical approaches that are for the most part semantically equivalent, which may lead to additional challenges in hypothesis formalization. For instance, an ANOVA with repeated measures and a linear mixed effects model are similar in intent but require two different function calls, one without a formula (e.g., `AnovaRM` in `statsmodels` [T14]) and another with (e.g., `mixedlm` in `statsmodels` [T14]). Even when considering only ANOVA, tools may provide similar syntax but implement different sums of squares procedures for partitioning variance (i.e., Type I, Type II, or Type III).⁶ By default, R's `stats` core package [T10] uses Type I, `statsmodels` [T14] uses Type II, and `SPSS` [T16] uses Type III. The three different sum of squares procedures lead to different F-statistics and p-values, which may lead analysts to different conclusions. More importantly, the procedures encode different conceptual hypotheses. If analysts have theoretical knowledge or conceptual hypotheses about the order of independent variables, tools defaulting to Type I (e.g., R's `stats` core library) align the model implementation with the conceptual hypotheses. However, if analysts do not have such conceptual hypotheses, tools' default behavior would execute (without error) and silently respond to a conceptual hypothesis different from the one the analyst seeks to test. In this way, syntactic and semantic mismatches can create a rift between model implementations and conceptual hypotheses. Furthermore, the impact of tools' "invisible" model implementation choices reinforces the interplay between conceptual and model implementation concerns during hypothesis formalization.

Theme: Computational issues.

Tools provide end-users with options for optimizers and solvers used to fit statistical models [T1, T2, T4, T6, T7, T8, T10, T11, T13, T16, T18], convergence criteria used for fitting models [T3, T6, T16, T18], and memory and CPU allocation [T2, T5, T12, T15], among more specific customizations. For instance, `lme4` [T6] allows analysts to specify the nonlinear optimizer and its settings (e.g., the number of iterations, convergence criteria, etc.) used to fit models. In `brms` [T2], analysts can also specify the number of CPUs to dedicate to fitting their models. Some computational settings are akin to performance optimizations, affecting computer utilization but not the results. However, not all computational changes are so well-isolated.

For example, the failure of a model's inference algorithm to converge (in **Model Implementation**) may prompt mathematical re-formulation (**Mathematical Equation**), which may cast **Observations about Data** in a new light, prompting **Causal Model** and **Conceptual Hypothesis** revision. In other words, computational failures and decisions may bubble up to conceptual

⁶Type I is (a) sensitive to the order in which independent variables are specified because it assigns variance sequentially and (b) allows interaction terms. Type II (a) does not assign variance sequentially and (b) does not allow interaction terms. Type III (a) does not assign variance sequentially and (b) allows interaction terms. For an easy-to-understand blog post, see [Kor19].

hypothesis revision and refinement, which may then trickle back down to model implementation iteration, and so on. In this way, computational control can be another entry into the dual-search process of hypothesis formalization.

In theory this low-level control could help analysts formalize nuanced conceptual hypotheses in diverse computational environments. However, we found that tools do not currently provide feedback on the ramifications of these computational changes, introducing a gulf of evaluation [Nor86]. Analysts can easily change parameters to fine-tune their computational settings, but how they should interpret their model implementations and revisions conceptually is unaddressed, suggesting opportunities for future tools to bridge the conceptual and model implementation gap.

4.4.3 Takeaways from the Analysis of Tools

Taken together, our analysis shows that tools can support a wide range of statistical models but expect analysts to have more statistical expertise than may be realistic. They provide limited guidance for analysts (i) to express and translate their conceptual and partially-formalized concerns and (ii) identify reasonable models. Tools also provide little-to-no feedback on the conceptual ramifications of model implementation iterations. These gaps reveal a misalignment between analysts' hypothesis formalization processes and tools' expectations and design. Possible reasons for this mismatch may be that tools do not scaffold or embody the dual-search nature of hypothesis formalization or leverage all the intermediate artifacts analysts may create (e.g., refined conceptual hypotheses, causal models, data observations, partial specifications, etc.) throughout the process.

4.5 Design Implications for Statistical Analysis Software

Our findings suggest three opportunities for tools to facilitate the dual-search process and align conceptual hypotheses with statistical model implementations at various stages of hypothesis formalization.

Meta-libraries: Connecting Model Implementations with Mathematical Equations

Specialized tools, although necessary for sophisticated statistical computation, require a steep learning curve. *Meta-libraries* could allow analysts to specify their statistical models in high-level code; execute the statistical models using the appropriate libraries in their knowledge bases; and then output library information, functions invoked, any computational settings used, the mathematical model that is approximated, and the statistical results. Libraries such as Parsnip [KVR20] have begun to provide a unified higher-level interface that allows analysts to specify a statistical model

using more “generically” named functions, parameter names, and symbolic formulae (when necessary). Parsnip then compiles and invokes various library-specific functions for the same statistical model.

Probabilistic programming languages (PPLs), such as Pyro [BCJ⁺19], Stan [CGDH⁺17], BUGS [LTBS00], PyMC [SWF16], already enable the development of meta-libraries. PPLs support modular specification of data, probabilistic models, and probabilistic hypotheses. Existing libraries, including `brms`, provide higher-level APIs whose syntax uses symbolic formulae, for instance, and compile to programs in a PPL (i.e., Stan in the case of `brms`).

As already seen in Parsnip and tools using PPLs, meta-libraries could bring three benefits. First, they would provide simpler, less fragmented interfaces to analysts while continuing to take advantage of tool specialization. Second, meta-libraries that output complete mathematical representations would more tightly couple mathematical representations with implementations, providing an on-ramp for analysts to expand their statistical knowledge. Third, meta-libraries that show the mathematical representations alongside underlying libraries’ function calls could show syntactical variation in underlying libraries, indirectly teaching analysts how they might express their statistical models in other tools, familiarizing analysts with new tools and statistical models, and even mend fragmented views of identical statistical approaches (e.g., ANOVA and regression).

Future meta-libraries could consider providing a higher-level, declarative interface that does not require analysts to write symbolic formulae. Designing such declarative meta-libraries would require formative elicitation studies (similar to natural programming studies such as [VAK18]) on declarative primitives that are memorable, distinguishable, and reliably understood. An additional challenge would lie in maintaining support for various libraries executed under the hood, especially as libraries change their APIs, which would strengthen the case for meta-libraries. Although meta-libraries would not solve the problems involved in understanding how computational settings affect statistical model execution or conceptual hypotheses, they could nevertheless provide scaffolding for analysts to more closely examine specific libraries, especially if multiple libraries execute the same statistical model but do not all encounter the same computational bottlenecks.

High-level Libraries: Expressing Conceptual Hypotheses to Bootstrap Statistical Model Implementations

The absence of tools for directly expressing conceptual hypotheses may be an explanation for why data workers in the lab study dove into statistical model implementation details. High-level libraries could allow analysts to specify data collection design (e.g., independent variables, dependent variables, controlled effects, possible random effects); variable data types; expected or known covariance relationships based on domain expertise; and hypothesized findings in a library-specific

grammar. High-level libraries could compile these conceptual and data declarations into weighted constraints that represent the applicability of various statistical approaches, in a fashion similar to Tea [JDR⁺19], a domain-specific language for automatically selecting appropriate statistical analyses for common hypothesis tests. Libraries could then execute the appropriate statistical approaches, possibly by using a meta-library as described above.

In addition to questions of how to represent a robust taxonomy of statistical approaches computationally, another key challenge for developing high-level libraries is identifying a set of minimal yet complete primitives that are useful and usable for analysts to express information that is usually expressed at different levels of abstraction: conceptual hypotheses, study designs, and possibly even partial statistical model specifications. For instance, even if a conceptual hypothesis is expressible in a library, it may be impossible to answer with a study design or partial statistical model that is expressed in the same program. An approach may be to draw upon and integrate aspects from existing high-level libraries and systems that aim to address separate steps of the hypothesis formalization process, such as Touchstone2 [EWBLM19] for study design and Tea and Statsplorer [WSVB15] for statistical analysis.

Bidirectional Conceptual Modeling: Co-authoring Conceptual Models and Statistical Model Implementations

Conceptual, or causal, modeling was difficult for the analysts in the lab study. Some even resisted conceptual modeling for fear of biasing their analyses. Yet, implicit conceptual models were evident in the hypotheses analysts chose to implement and the sub-hypotheses researchers articulated in the content analysis.

Mixed-initiative systems that make explicit the connection between conceptual models and statistical model implementations could facilitate hypothesis formalization from either search process and allow analysts to reflect on their analyses without fear of bias. For example, a mixed-initiative programming environment could allow analysts to write an analysis script, detect data variables in the analysis scripts, identify how groups of variables co-occur in statistical models, and then visualize conceptual models as graphs where the nodes represent variables and the edges represent relationships. The automatically generated conceptual models would serve as templates that analysts could then manipulate and update to better reflect their internal conceptual models by specifying the kind of relationship between variables (e.g., correlation, linear model, etc.) and assigning any statistical model roles (e.g., independent variable, dependent variable). As analysts update the visual conceptual models, they could evaluate script changes the system proposes. In this way, analysts could externally represent their causal models while authoring analysis scripts and vice versa.

Although bidirectional programming environments already exist for vector graphics creation [HLC19], they have yet to be realized in mainstream data analysis tools. To realize bidirectional, automatic conceptual modeling, researchers would need to address important questions about (i) the visual grammar, which would likely borrow heavily from the causal modeling literature; (ii) program analysis techniques for identifying variables and defining co-occurrences (e.g., line-based vs. function-based) in a way that generalizes to multiple statistical libraries; and (iii) adoption, as analysts who may benefit most from such tools (likely domain non-experts) may be the most resistant to tools that limit the number of “insights” they take away from an analysis.

4.6 Discussion

Hypothesis formalization is a dual-search process of translating conceptual hypotheses into statistical model implementations. Due to constraints imposed by domain expertise, data, and tool familiarity, the same conceptual hypothesis may be formalized into different model implementations. A single model implementation may be useful for making multiple statistical inferences. The same model implementation may also formalize two possibly opposing hypotheses. To navigate these constraints, analysts use problem-solving strategies characteristic of the larger scientific discovery process [KD88; SK95]. As such, hypothesis formalization exemplifies how data science is a design practice.

At a conceptual level, hypothesis formalization involves *hypothesis refinement*, which, to use Schunn and Klahr’s language [SK95], is a *scoping* process. In the formative content analysis, we found that researchers *decomposed* their research goals and conceptual hypotheses into specific, testable sub-hypotheses and *concretized* constructs using proxies, born of theory or available data. Also, we found that analysts in the lab study also quickly converged on the need to specify established proxies or develop them based on the data schema presented. In hypothesis formalization, scoping incorporates domain- and data-specific observations to qualify the conceptual scope of researchers’ hypotheses. In other words, hypothesis refinement is an instance of *means-end analysis* [NS⁺72], a problem-solving strategy that aims to recursively change the current state of a problem into sub-goals (i.e., increasingly specific objectives) in order to apply a technique (i.e., a particular statistical model) to solve the problem (i.e., test a hypothesis).

At the other computational endpoint of hypothesis formalization, *statistical model implementation* also involves iteration. Through our analysis of software tools, we found that analysts must not only select tools among an array of specialized and general choices but also navigate tool-specific taxonomies of statistical approaches. These tool taxonomies may both differ from and inform analysts’ personal categorizations, potentially explaining why analysts in our lab study relied on their

personal taxonomies and tools. Based on their prior experience, analysts engage in *analogical reasoning* [HHNT89], finding parallels between the present analysis problem’s structure and previously encountered ones or ones that fit a tool’s design easily.

Upon selecting a statistical function, analysts may tune computational settings, choose different statistical functions or approaches, which they may tune, and so on. In this way, the model implementation loop in hypothesis formalization captures the “debugging cycles” analysts encounter, such as the census researcher in the introduction. The tool ecosystem as a whole supports diverse model implementations, even for the same mathematical equation. However, the tool interfaces provide low-level abstractions, such as interfaces using mathematical formulae that, based on our observations in the lab study, do not support the kind of higher-level conceptual reasoning required of hypothesis formalization.

4.7 Future Work

The steps, considerations, and strategies we have identified are domain-general. Domain-specific expertise likely influences how quickly analysts switch between steps and strategies during the dual-search process. Domain experts, including researchers in our content analysis, may know which statistical model implementations and computational settings to use *a priori* and design their studies or specify their conceptual hypotheses in light of these expectations — incorporating means-end analysis and analogical reasoning strategies — more quickly. It may be these insights that analysts in our lab study sought when they looked online for conceptual and statistical help.

Future work could observe how domain experts perform hypothesis formalization and characterize when and how analysts draw upon their own or collaborators’ expertise to circumvent iterations or justify early scoping decisions. These insights may also shed light on how pre-registration expectations and practices could be made more effective. Given the level of detail required of some pre-registration policies, researchers likely engage in a version of the hypothesis formalization process we have identified prior to registering their studies. Knowing how pre-registration fits into the hypothesis formalization process could improve the design and adoption of pre-registration practices.

Future work could also explore how hypothesis formalization may differ in machine learning settings. In this chapter, our focus was on how analysts answer domain questions and test hypotheses using statistical methods and their domain knowledge. Our findings may not generalize to settings or methods where domain knowledge is less important, such as deep learning and other machine learning-based approaches.

Finally, our findings suggest opportunities for future tools to bridge steps involved in hypothesis formalization and guide analysts towards reasonable model implementations. Our analysis of

tools suggest possibilities for tools to connect model implementations to their mathematical representations through meta-libraries, provide higher-level abstractions for more directly expressing conceptual hypotheses, and support automated conceptual modeling. Future system development and user testing are necessary to validate these implications and more readily support analysts translate their conceptual hypotheses into statistical model implementations.

4.8 Summary of Contributions

The empirical studies that led us to articulate the theory of hypothesis formalization illustrates the key challenge to authoring data analyses: Analysts must translate their implicit domain knowledge into statistical specifications that they can implement and execute in code. As we saw in the lab study, analysts often resort to changing their hypotheses or research questions to what they can implement or get stuck on how to represent their conceptual knowledge in statistical models, highlighting the dual-search nature of hypothesis formalization. Furthermore, the summary of hypothesis formalization (i.e., Figure 4.1) serves as a device for (i) interpretation—to explain where and how analysts struggle in authoring statistical analyses—and (ii) inspiration—to inspire new approaches and systems to authoring data analyses.

Our theory of hypothesis formalization highlights the discrepancy between analysts' goals and the statistical software tools available to them. While analysts want to understand their data to better understand their domains or make decisions, the current ecosystem prioritizes mathematical expressivity and computational control, features that are likely desirable for statistical experts but not novices.

As a result, designing new data analysis tools to gather conceptual knowledge and translate them into statistical analyses is a promising approach for statistical non-experts. In this way, hypothesis formalization retrospectively validates our design in Tea, where its constraint-based runtime system provided automated reasoning for Null Hypothesis Significance Tests. In order to support more complex research questions, additional methods of explicitly grappling with more conceptual knowledge and reasoning about different classes of statistical analyses is necessary. We tackle this challenge for generalized linear models with and without mixed effects in Tisane.

This work was in collaboration with Nicole de Moura, Melissa Birchfield, Jeffrey Heer, and René Just. It was originally published in ACM Transactions of Computer-Human Interaction (TOCHI) 2022 [JBDM⁺ 22a] and presented at ACM CHI 2022.

Chapter 5

Tisane: Authoring statistical models via formal reasoning from conceptual and data relationships

Authoring statistical models requires analysts to jointly reason about their conceptual domain knowledge, statistical methods, and analysis implementations in code, as our theory of hypothesis formalization describes. For instance, scientists carefully consider which covariates to include in statistical models based on their prior knowledge of confounding. However, analysts' conceptual knowledge is often kept implicit. Analysts gravitate towards statistical specifications they are familiar with, even if the analyses are sub-optimal or do not assess their hypotheses, as we saw in the previous chapter. Finally, ease of implementation further constrains the statistical models that analysts try and use. These issues are especially salient for domain experts who lack deep statistical or programming expertise (e.g., many researchers).

Existing statistical software exacerbate these issues because they do not allow analysts to externalize their implicit conceptual knowledge, receive guidance on analysis approaches, and help authoring low-level statistical modeling code Section 4.4. Our work on Tisane hypothesizes that in order to address these issues, software tools should capture analysts' implicit conceptual models and use them to derive statistical models.

*Conceptual models*¹ are often-informal representations of variable relationships (e.g., list of variable relationships, process diagrams, graphs), describing the underlying data generating process. Conceptual models are difficult to reason about during statistical analysis. Their implications on statistical modeling are not obvious, especially to statistical non-experts. For example, the impact

¹Richard McElreath calls these implicit assumptions *process models* [McE20]. We use the term *conceptual models* in order to contrast from statistical models.

of conceptual assumptions may only become apparent after fitting multiple statistical models, if at all. Without explicitly grappling with conceptual models prior to authoring statistical models, analysts run the risk of introducing inconsistencies between their domain knowledge and statistical models, which can lead to unintentionally answering a different research question and asserting a conceptual model based on preferred results (i.e., HARKing [1]).

To facilitate more accurate hypothesis formalization and analysis, we asked, **How might we derive (initial) statistical models from conceptual models?**. Inferring a statistical model raises two technical challenges: (1) How do we elicit the information necessary for inferring a statistical model? and (2) How do we infer a statistical model, given this information? We explore and address these issues by iteratively designing, developing, and evaluating **Tisane, a system for implementing generalized linear models (GLMs) and generalized linear mixed-effects models (GLMMs) from explicit statements of implicit conceptual assumptions**.

The first implementation of Tisane (Section 5.3) was as an open-source Python package available on pip (see [JS23]). Case studies and real-world use of Tisane demonstrated not only the viability but also the desirability of tool support for authoring statistical models from conceptual models (Section 5.4). Therefore, we explored how to further improve Tisane’s programming and interaction models to better suit novice analysts (Section 5.6) and released a second version in R as the rTisane library [JMHJ23] (see [JS23]). The R implementation allowed us to more directly compare rTisane to a scaffolded workflow using widely used linear modeling libraries, including the `lme4`, in R.

Tisane provides a **study design specification language** for expressing relationships between variables. Tisane compiles the explicitly stated relationships into an internal **graph representation** and then traverses the graph to infer candidate GLMs/GLMMs based on recommendations from the graphical causal reasoning community. Analysts can then query Tisane for a statistical model that explains a specific dependent variable from a independent variable of interest. Tisane helps analysts disambiguate their input conceptual models and an output statistical model script for fitting a valid GLM/GLMM. In this way, Tisane focuses analysts on reflecting on and externalizing their implicit conceptual assumptions and checks that analysts do not overlook relevant variables, such as potential confounders or data clustering, that could compromise generalizability of statistical results.

5.1 Background and Related work

At the heart of Tisane is the goal to derive statistical models from conceptual models. To do so, Tisane relies on transforming aspects of analysts’ expressed conceptual models into causal graphs. There are multiple frameworks for reasoning about causality [Rub04; Pea95a]. One widespread

approach is to use directed acyclic graphs (DAGs) to encode conditional dependencies between variables [Pea95b; GPR99; Spi94; SRM⁺96]. If analysts can specify a formal causal graph, Pearl’s “backdoor path criterion” [Pea95a; P⁺00] explains the set of variables that control for confounding. However, in practice, specifying proper causal DAGs is challenging and error-prone for domain experts who are not also experts in causal analysis [SSY20]. Empirical findings may be inconclusive or ambiguous in the causal relationships they suggest [SV18]. Statistical non-experts also lack guidance on which variables and relationships to include [VDN⁺13]. Despite having important domain knowledge, analysts do not have interfaces that allow them to express what they know in a way that is approachable to them. Therefore, Tisane does not expect analysts to specify a formal causal graph. Instead, analysts can express causal relationships as well as more ambiguous relationships between variables in the study design specification language.

Furthermore, prior work in the causal reasoning literature shows how linear models can be derived from causal graphs to make statistical inferences and test the motivating causal graph [SRM⁺96; Spi94]. Recently, VanderWeele proposed the “modified disjunctive cause criterion” [Van19] as a new heuristic for researchers without a clearly accepted formal causal model to identify confounders to include in a linear model, for example. The criterion identifies confounders in a graph based on expressed causal relationships. The first release of Tisane (Section 5.3) applies the modified disjunctive cause criterion when suggesting variables to include in a GLM or GLMM. Tisane does not automatically include variables to the statistical models because substantive domain knowledge is necessary to resolve issues of temporal dependence between variables, among other considerations [Van19]. In the second release of Tisane (Section 5.7), we use the more recent recommendations from Cinelli, Forney, and Pearl [CFP20] for controls in regression models. To guide analysts through the suggestions, Tisane provides analysts with explanations to aid their decision making during disambiguation.

Importantly, generalized linear models with or without mixed effects are not formal causal analyses. Tisane does not calculate average causal effect or other causal estimands. Rather, Tisane only utilizes insights about the connection between causal DAGs and linear models to guide analysts towards including potentially relevant confounders in their GLMs grounded in domain knowledge.

5.1.1 Statistical scope

Generalized linear models (GLMs) and generalized linear models with mixed effects (GLMMs) are meaningful targets because they are commonly used (e.g., in psychology [LA15; CCWA13], social science [KKdL98], and medicine [BBC⁺09; BLST13]) yet are easy to misspecify for statistical experts and non-experts alike [BLST13; CCWA13]. We designed Tisane to support researchers who are domain experts capable of supplying conceptual and data collection information but lack the

statistical expertise or confidence to author GLM/GLMMs accurately. Both GLMs and GLMMs consist of (i) a *model effects structure*, which can include main and interaction effects and (ii) *family* and *link* functions. The family function describes how the residuals of a model are distributed. The link function transforms the predicted values of the dependent variable. This allows modeling of linear and non-linear relationships between the dependent variable and the predictors. In contrast to transformations applied directly to the dependent variable, a link function does not affect the error distributions around the predicted values. The key difference between GLMs and GLMMs is that GLMMs contain random effects in their model effects structure. Random effects describe how individuals (e.g., a study participant) vary and are necessary in the presence of hierarchies, repeated measures, and non-nesting composition (5.3.1)².

Both GLMs and GLMMs assume that (i) the variables involved are linearly related, (ii) there are no extreme outliers, and (iii) the family and link functions are correctly specified. In addition, GLMs also assume that (iv) the observations are independent. Tisane’s interactive compilation process guides users through specifying model effects structures, family and link functions to satisfy assumption (iii), and random effects only when necessary to pick between GLMs and GLMMs and satisfy assumption (iv).

We scoped Tisane to GLM and GLMMs because they encompass a large scope of statistical models such that our research contributions are widely applicable and substantial. In addition, given that GLMs and GLMMs can represent common Null Hypothesis Significance Tests (in Tea), Tisane generalizes our approach in Tea. Tisane gives further evidence of the benefit of conceptual programming abstractions and automated reasoning for authoring statistical analyses.

5.2 Early Design Process

Tisane’s first released DSL was the result of an iterative design process, including informal usability critiques of language design and a user study with three researchers. We describe the process further below.

With Tisane’s graph specification language, we aimed to collect the necessary information to infer a GLM/GLMM and to provide a straightforward way of collecting it. We consulted statistical best practices on how to construct valid GLMs [KKdL98; Bar13; BLST13; McE20], which led us to two sets of variable relationships: *conceptual relationships*, specifically about causal and correlational relationships to explain using a GLM/GLMM and *data measurement relationships* about the frequency of observations per observational unit (or “level”) and how observations may be clustered

²Traditionally, the term “mixed effects” refers to the simultaneous presence of “fixed” and “random” effects in a single model. We try to avoid these terms as there are many contradictory usages and definitions [Gel05]. When we do use these terms, we use the definitions from Kreft and De Leeuw [KKdL98].

(e.g., nesting).

We conducted an exploratory survey of 12 study design and data collection packages. We identified these libraries using word of mouth and bibliographic references. Eight libraries focused on the controlling the presentation of stimuli and trials (lower-level). Five were focused on the distribution of conditions (e.g., within-subjects vs. between-subjects) and frequency of measures (higher-level). We prototyped Tisane’s SDSL based on the constructs common across these tools.

To better understand how using variable relationships to author statistical models affects data analysis workflows, we tested an earlier prototype of Tisane with three computer science researchers (in AI, HCI, and systems, whom we refer to as P1, P2, and P3, respectively). We were concerned that we were redistributing the difficulty of authoring GLMs/GLMMs from specifying them directly to expressing potentially obscure variable relationships.

All three researchers reported that the study design specification language was straightforward. P2 remarked, “The API is very simple and elegant. It’s very intuitive. It gets me really thinking about what’s the essential or most important part of the analysis.” Needing to explicitly state variable relationships in Tisane prompted P1 and P2 to think more critically about their domain and discover new analysis paths (P1, P2, P3). For example, Tisane helped P1, who previously had erroneously believed multiple t-tests with Bonferroni corrections were more appropriate than a GLM for his data, realize how a GLM could have helped him answer questions he had not had the foresight to ask beforehand. We were encouraged to see researchers reap additional benefits of having to specify variable relationships.

Earlier versions of Tisane had a more extensive API that distinguished between observations and experimental treatments and provided multiple ways to specify the same types of relationships. We observed that the researchers gravitated toward a smaller subset of language constructs around unit and measure declaration, so we introduced explicit types for units and measures and removed redundant functions.

5.3 First Release

Tisane provides a *study design specification language (SDSL)* for expressing relationships between variables. There are two key challenges in designing a specification from which to infer statistical models: (1) determining the set of relationships that are essential for statistical modeling and (2) determining the level of granularity to express relationships.

Table 5.1: Overview of study design tools that informed Tisane’s study design specification language. The first five tools provide higher-level abstractions. They are designed to help researchers reason about their study designs more holistically. The latter eight tools are lower-level and are more focused on stimuli, trials, and progressions between trials. *JsPsych is the base package to which JsPsychR, xprmtnr, and Jaysire provide wrappers and extensions.

| Tool | Support provided |
|--|---|
| Edibble [Tan21] JMP Design of Experiments [SAS20b] Gosset [SH17] DeclareDesign [BCCH19b] Touchstone2 [EWBLM19] | reason about end-to-end experimental design, create data collection schema use templates for experiments, some design optimization, some help with modeling search for optimal study design simulate data, specify and reason about designs statistically design controlled experiments while reasoning about randomization and statistical power |
| Formr [AWT20] psychTestR [Har20] Psychopy [Pei07] | design online survey questions and flow create trials, specify “timelines” for how trials should progress control how (visual) stimuli are presented, trials, and trial progression in an online experiment |
| Psychtoolbox [BSG12] JsPsych* [DL15b] JsPsychR [Cru19] xprmtnr [Nav19] Jaysire [Nav21] | control stimuli in an online experiment, especially for neuroscience create and control trials and stimuli for online experiments |

5.3.1 Study design specification language and graph representation

In Tisane’s SDSL, analysts can express conceptual and data measurement relationships between variables. Both are necessary to specify the domain knowledge and study designs from which Tisane infers statistical models.

Variables

There are three types of data variables in Tisane’s SDSL: (i) units, (ii) measures, and (iii) study environment settings. The **Unit** type represents entities that are observed and/or receive experimental treatments. In the experimental design literature, these entities are referred to as “observational units” and “experimental units,” respectively. Entities can be both observational and experimental units simultaneously, so the SDSL does not provide more granular unit sub-types. The **Measure** type represents attributes of units and must be constructed through their units, e.g., `age = adult.numeric('age')`. Measures are proxies (e.g., minutes ran on a treadmill) of underlying constructs (e.g., endurance). Measures can have one of the following data types: numeric, nominal, or ordinal. Numeric measures have values that lie on an interval or ratio scale (e.g., age, minutes ran on a treadmill). Nominal measures are categorical variables without an ordering (e.g., race). Ordinal measures are ordered categorical variables (e.g., grade level in school). We included these data types because they are commonly taught and used in data analysis. The **SetUp** type represents study environment settings that are neither units nor measures. For example, time is often

an environmental variable that differentiates repeated measures but is neither a unit nor a measure of a specific unit.

Relationships between Variables

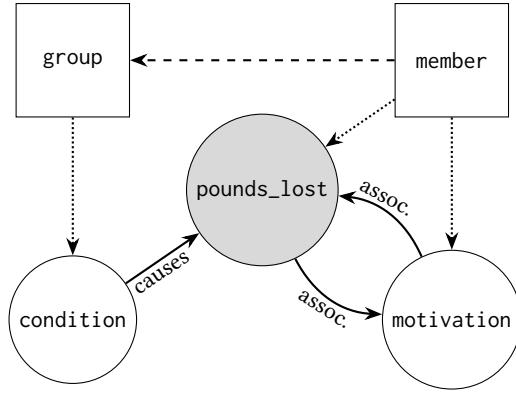


Figure 5.1: The graph representation of the variables and relationships from the usage scenario. causes edges are labeled with “causes”. associates_with edges are labeled with “assoc.” Dashed edges indicate nests_within relationships, and dotted edges indicate has relationships.

In Tisane’s SDSL, variables have relationships that fall into two broad categories: (1) *conceptual relationships* that describe how variables relate theoretically and (2) *data measurement relationships* that describe how the data was, or will be, collected. Below, we define each of the relationships in Tisane’ SDSL and describe how Tisane internally represents these relationships as a graph (as illustrated in 5.3.1). 5.1 shows the graph representation constructed from the usage scenario.

Tisane’s graph IR is a directed multigraph. Nodes represent variables, and directed edges represent relationships between variables. Tisane internally uses a graph intermediate representation (IR) because graphs are widely used for both conceptual modeling and statistical analysis, two sets of considerations that Tisane unifies.

Tisane’s graph IR differs from two types of graphs used in data analysis: causal DAGs and path analysis diagrams. Unlike causal DAGs, Tisane’s graph IR allows for non-causal relationships, moderating relationships (i.e., interaction effects), and data measurement relationships that are necessary for inferring random effects. Unlike path analysis diagrams that allow edges to point to other edges to represent interaction effects, Tisane represents interactions as separate nodes and only allows nodes as endpoints for edges. These design decisions simplify our statistical model inference algorithms and their implementation.

Conceptual relationships. Tisane’s SDSL supports three conceptual relationships: causes, associates with, and moderates. Analysts can express that a variable **causes** or is **associated with**

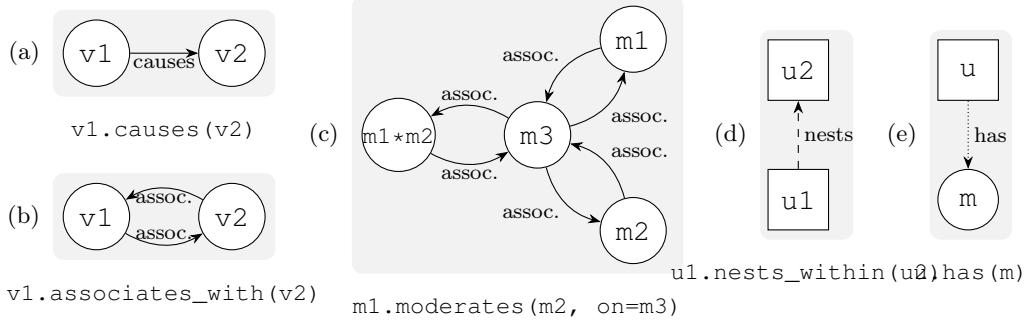
(but not directly causally related to) another variable. Variables associated with the dependent variable, for example, may help explain the dependent variable even if the causal mechanism is unknown. If analysts are aware of or suspect a causal relationship, they should use `causes`.

We chose to support both causal and associative relationships because formal causal DAGs are difficult for domain experts to specify [SSY20; SV18; VDN⁺13], prior work has observed that researchers already use informal graphs that contain associative relationships when reasoning about their hypotheses and analyses [JBDM⁺22b], and GLMs/GLMMs can represent non-causal relationships. Finally, analysts can also express interactions where one (or more) variable (the *moderating variables*) **moderates** the effect of a *moderated variable* on another variable (the *target variable*).

Mediation relationships (where one variable influences another through a middle variable) are another common conceptual relationship. Tisane does not provide a separate language construct for mediation because mediations are expressible using two or more causal relationships. Furthermore, mediation analyses require specific analyses, such as structural equation modeling [Hoy95], that are out of Tisane’s scope.

In the graph IR, a `causes` relationship introduces a causal edge from one node, the cause, to another node, the effect (5.3.1(a)). Because a variable cannot be both the cause and effect of the same variable, any pair of nodes can only have one causal edge between them. Furthermore, from a formal causal analysis perspective, associations may indicate the presence of a hidden, unobserved variable that mediates the causal effect of a variable on another or that influences two or more variables simultaneously. Thus, rather than inferring or requiring analysts to specify hidden variables, which may be unknown and/or unmeasurable, the `associates_with` relationship introduces two directed edges in opposing directions, representing the bidirectionality of association (5.3.1(b)). A `moderates` relationship creates a new node that is eventually transformed into an interaction term in the model, introduces associative edges between the new interaction node and the target (variable) node, creates associative edges between the moderated variable’s node and the target node, and adds associative edges between the moderating variables’ nodes and the target node if there is not a causal or associative edge already (5.3.1(c)). Furthermore, each interaction node inherits the attribution edges from the nodes of the moderating variables that comprise it. This means that every interaction node is also the attribute of at least one unit.³

³In statistical terms, this means that within-level interactions have one unit while cross-level interactions may have two or more units.



Data measurement relationships. Study designs may have clusters of observations that need to be modeled explicitly for external validity. For example, in a within-subjects experiment, participants provide multiple observations for different conditions. An individual's observations may cluster together due to a hidden latent variable. Such clustering may be imperceptible during exploratory data visualization of a sample but can threaten external validity. GLMMs can mitigate three common sources of clustering that arise during data collection [GH06; KKdL98; Coh88]:

- **Hierarchies** arise when one observational/experimental unit (e.g., adult) nests within another observational/experimental unit (e.g., group). This means that each instance of the nested unit belongs to one and only one nesting unit (many-to-one).
- **Repeated measures** introduce clustering of observations from the same unit instance (e.g., participant).
- **Non-nesting composition** arises when overlapping attributes (e.g., stimuli, condition) describe the same observational/experimental unit (e.g., participant) [GH06].

The above sources of clustering pose three problems for analysts. First, analysts must have significant statistical expertise to identify when data observations cluster. Second, they must know how to mitigate these clusters in their models. Third, with this knowledge, analysts must figure out how to express these types of clustering in their analytical tools. Even if analysts are not able to identify clustered observations, they are knowledgeable about how data were collected.

Thus, Tisane addresses the three problems by (i) eliciting data measurement relationships from analysts to infer clusters and (ii) formulating the maximal random effects structure, optimizing for external validity (5.3.2). Below, we describe language features for expressing data measurement relationships.

Nesting relationships: Hierarchies **Hierarchies** arise when a unit (e.g., an adult) is nested within another unit (e.g., an exercise group). Researchers may collect data with hierarchies to study

individual and group dynamics together or as a side effect of recruitment strategies. To express such designs, Tisane provides the `nests_within` construct. Conceptually, nesting is strictly between observational/experimental units, so Tisane type checks that the variables that nest are both Units. In the graph IR, a nesting relationship is encoded as an edge between two unit nodes (5.3.1(d)). There is one edge from the nested unit (e.g., `adult`) to the nesting unit (e.g., `group`) ⁴.

Frequency of measures: Repeated measures, Non-nesting composition When a measure is declared through a unit, Tisane adds an attribution edge (“has”) from a unit node to a measure node (5.3.1(e)). A unit’s measure can be taken one or more times in a study. The frequency of measurement is useful for detecting repeated measures and non-nesting composition. In **repeated measures** study designs, each unit provides multiple values of a measure, which are distinguished by another variable, usually time. **Non-nesting** [GH06] composition arises when measures describing the same unit overlap. For example, HCI researchers studying input devices might design them to utilize different senses (e.g., touch, sight, sound). Participants in the study may be exposed to multiple different devices, which act as experimental conditions of senses. The conditions are intrinsically tied to the devices, and participants can be described as having both conditions and devices, which overlap with one another. Such study designs introduce dependencies between observations [Cla73] and hence violate the assumption of independence that GLMs make.

When analysts declare Measures, they specify the frequency of the observation through the `number_of_instances` parameter. This parameter accepts an integer, variable, a Tisane `Exactly` operator, or a Tisane `AtMost` operator. By default, the parameter is set to one. The `Exactly` operator represents the exact number of times a unit has a measure. The `AtMost` operator represents the maximum number of times a unit has a measure. Both operators are useful for specifying that a measure’s frequency depends on another variable, which is expressible through the `per` function. For example, participants may use two devices *per* condition assigned: `device = subject.nominal('Input device', number_of_instances=ts.Exactly(2).per(condition))`. The `per` function uses the Tisane variable’s cardinality by default but can instead use a data variable’s `number_of_instances` by specifying `use_cardinality=False` as a parameter to `per`. Moreover, specifying a measure’s `number_of_instances` to be an integer is syntactic sugar for using the `Exactly` operator. Specifying a variable is syntactic sugar for expressing `ts.Exactly(1).per(variable)`.

To determine the presence of repeated measures or non-nesting composition, Tisane computes the `number_of_instances` of measures and their relationship to other measures. Measures that are declared with `number_of_instances` equal to one are considered to vary between-unit.

⁴The GitHub repo contains a gallery of examples that include nesting relationships.

Measures that are declared with `number_of_instances` greater than one or a variable with cardinality greater than one are considered to vary within-unit as repeated measures. If there are instances of a measure per another measure sharing the same unit, the measures are non-nesting.

5.3.2 Statistical model inference: Interactively querying the graph IR

After specifying variable relationships, analysts can query Tisane for a statistical model. Queries are constructed by specifying a study design with a dependent variable (the value to be predicted) and a set of independent variables (predictors). Tisane processes the query and generates a statistical model in four phases: (1) preliminary conceptual checks that validate the study design, (2) inference of possible effects structures and family and link functions, (3) input elicitation to disambiguate possible models, and (4) generation of a final executable script, and a record of decisions during disambiguation. Given that the interactive process begins with an input program using Tisane and outputs a script for fitting a GLM or GLMM, we call this process *interactive compilation*.

Preliminary checks

At the beginning of processing a query, Tisane checks that every input study design is well-formed. This involves two conceptual correctness checks. First, every independent variable (IV) in the study design must either cause or be associated with the dependent variable (DV) directly or transitively. Second, the DV must not cause any of the IVs, since it would be conceptually invalid to explain a cause from any of its effects. If any of the above checks fail, Tisane issues a warning and halts execution. By using these two checks, the Tisane compiler avoids technically correct statistical models that have little to no conceptual grounding (*DG1 - Conceptual knowledge*). If the checks pass, Tisane proceeds to the next phase.

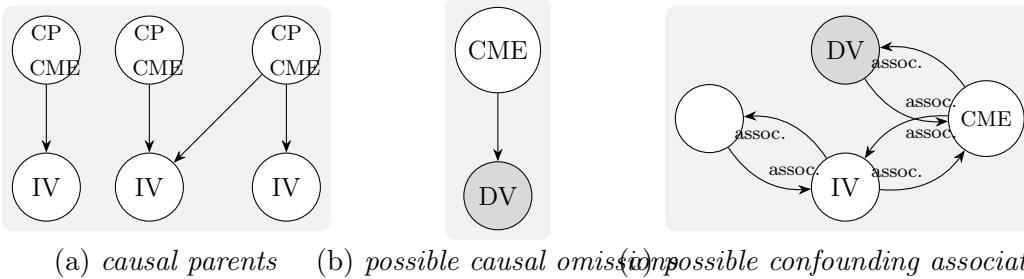
Candidate statistical model generation

A GLM/GLMM is comprised of a model effects structure, family function, and link function. The model effects structure may consist of main, interaction, and random effects. Tisane utilizes variables' conceptual relationships to infer candidate main and interaction effects and data measurement relationships to infer random effects. Tisane infers family and link functions based on the data type of the DV in the query. The candidate statistical models that Tisane generates, based on the graph and query, seed an interactive disambiguation process.

The purpose of identifying candidate main effects beyond the ones analysts may have specified is to provoke consideration of erroneously omitted variables that are conceptually relevant and pre-empt potential confounding and multicollinearity issues that may arise.

Deriving Candidate Main Effects In a query to infer a statistical model, analysts specify a single dependent variable and a set of one or more IVs. After passing the checks described in 5.3.2, the query’s independent variables are considered candidates. In addition, Tisane derives three additional sets of candidate main effects intended to control for confounding variables in the output statistical model⁵. The first two sets below are from the “modified disjunctive cause criterion” [Van19]:

- **Causal parents.** For each IV in the query, Tisane finds its causal parents (see 5.3.2(a)).
- **Possible causal omissions.** Tisane looks to see if any other variables not included as IVs cause the DV (see in 5.3.2(b)). They are relevant to the DV but may have been erroneously omitted.
- **Possible confounding associations.** For each IV, Tisane looks for variables that are associated with both the IV and the DV (see in 5.3.2(c)). Because associations between variables can have multiple underlying causal structures, Tisane recommends variables with associative relationships with caution. Tisane issues a warning describing when not to include such a variable in the GUI.



Using the above rules, Tisane suggests a set of variables that are likely confounders of the variables of interest expressed in the query. There may be additional confounders due to unmeasured or unexpressed variables that are either not known or excluded from the graph. Tisane never automatically includes the candidate main effects in the output statistical model. Analysts must always specify a variable as an IV in the query or accept a suggestion (*DG3 - Guidance and control*).

If a graph only contains associates edges then the candidate main effects Tisane suggests are those that are directly associated with both the DV and an IV. If a graph has only causal edges, Tisane would suggest variables that directly cause the DV but were omitted from the query and the causal parents of IVs in case the parents exert causal influence on the DV through the IV or another variable that is not specified.

⁵Tisane currently treats each input IV as a separate “exposure” variable for which to identify confounders. Tisane then combines all confounders into one statistical model.

The total set of main effects, including variables the analyst has specified as IVs in their query and candidate main effects, are used to derive candidate interaction effects and random effects, which we discuss next.

Deriving Candidate Interaction Effects An interaction between variables means that the effect of one variable (the *moderated* variable) on a *target* variable is moderated by another (non-empty) set of variables (the *moderating* variables). Tisane’s SDSL already provides a primitive, `moderates`, to express interactions. As such, Tisane’s goal in suggesting candidate interaction effects is to help analysts avoid omissions of conceptual relationships that are pertinent to an analyst’s research questions or hypotheses (*DG1 - Conceptual knowledge*). Candidate interaction effects are the interaction nodes whose (i) moderated and moderating variables include two or more candidate main effects and (ii) target variable is the query’s DV.

Deriving Candidate Random Effects Random effects occur when there are clusters in the data, which occur when we have repeated measures, nested hierarchies, or non-nesting composition (as defined in subsection 5.3.1). Tisane implements Barr et al.’s recommendations for specifying the maximal random effects structure of linear mixed effects models for increasing the generalizability of statistical results [BLST13; Bar13].

To derive random effects, Tisane focuses on the data measurement edges in the graph IR. Using the graph IR, Tisane identifies unit nodes, looks for any nesting edges among them, and determines within- or between-subjects measures based on the frequency of observations for units. From these, Tisane generates random intercepts of units for the unit’s measures that are between-subjects as well as the unit’s measures that are within-subjects where each instance of the unit has only one observation per value of another variable. Tisane generates random slopes of a unit and its measure for all measures that are within-subjects where each instance of the unit has multiple observations per value of another variable. For interaction effects, random slopes are included for the largest subset of within-subjects variables (see [Bar13]). Tisane handles correlation of random slopes and intercepts during disambiguation (subsection 5.3.2). Maximal random effects may lead to model convergence issues that analysts address by later removing or adding independent variables and random effects. Nevertheless, starting with a maximal, valid model is important for ensuring that future revisions are also valid (*DG2 - Validity*).

Deriving Candidate Family and Link Functions The DV’s data type determines the set of candidate family and link functions. For example, numeric variables cannot have binomial or multinomial distributions. Similarly, nominal variables are not allowed to have Gaussian distributions.

Furthermore, each family has a set of possible link functions. For example, a Gaussian family distribution may have an Identity, Log, or Square Root link function. The statistics literature documents possible combinations of family and link functions for specific data types [NW72].

Tisane includes common family distributions as candidate families and their applicable link functions. In its current implementation, Tisane relies on `statsmodels` [SP10] for GLMs and `pymer4` [Jol18] for GLMMs. As such, Tisane is limited to the family and link function pairings implemented in these libraries. Table 5.2 lists the family and link functions these libraries currently supports. As `statsmodels`' and `pymer4`'s support for GLMs grows in the future, Tisane can be extended.

Table 5.2: The available family and link functions in Tisane. Tisane generates code to fit models using `statsmodels` and `pymer4`. The package `statsmodels` supports GLMs without mixed-effects and a wider variety of family and link function combinations. The package `pymer4` supports GLMs with mixed effects and has much more limited support for family and link functions. As `statsmodels` and `pymer4` add more support, Tisane can be extended.

| Family functions | Link functions (*default) | |
|-------------------|--|--|
| | Generalized linear models without mixed effects (<code>statsmodels</code>) | Generalized linear models with mixed effects (<code>pymer4</code>) |
| Gaussian | Identity*, Inverse, Log | Identity* |
| Inverse Gaussian | Identity, Inverse, Inverse Squared*, Log | Inverse Squared* |
| Gamma | Identity, Inverse*, Log | Inverse* |
| Poisson | Identity, Log*, Square Root | Log* |
| Binomial | Cauchy, CLogLog, Log, Logit*, Probit, | Logit* |
| Negative Binomial | Identity, Log*, Logit, Probit | N/A |
| Tweedie Family | Log*, Power | N/A |

Eliciting Analyst Input for Disambiguation

The disambiguation process provides an opportunity for analysts to explore the space of generated models based on their original query. Given our design considerations to prioritize conceptual knowledge (*DG1 - Conceptual knowledge*) and give analysts guidance (*DG3 - Guidance and control*), we designed a GUI to scaffold analysts' reasoning and elicit their input. For versatility, we implemented Tisane's GUI using Plotly Dash [Com23]. Analysts can either execute their Tisane programs and use the GUI inside a Jupyter notebook (no additional widgets needed) or run their Tisane programs in an IDE or terminal, in which case Tisane will open the GUI in a web browser.

Candidate statistical models are organized according to (i) independent variables (main effects and interaction effects), (ii) data clustering (random effects), and (iii) data distribution (family and link functions). In the main effects tab, Tisane asks analysts if they would like to include additional or substitute main effects that Tisane infers to be conceptually relevant. In the interaction effects tab, Tisane suggests moderating relationships to include but does not automatically include them

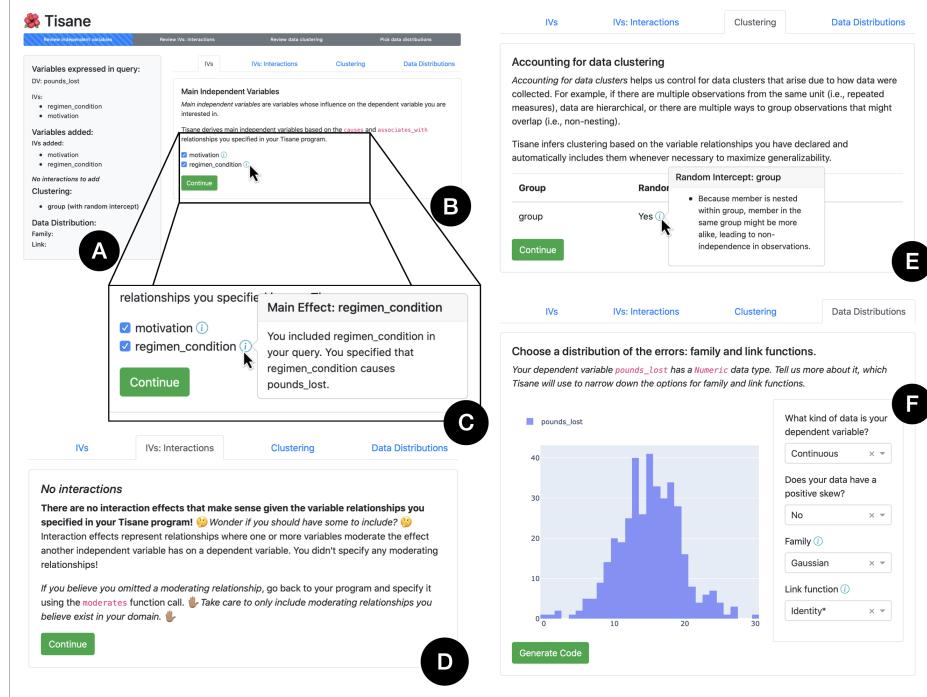


Figure 5.2: Example Tisane GUI for disambiguation. Tisane asks analysts disambiguating questions about variables that are conceptually relevant and that analysts may have overlooked in their query. (A) The left hand panel gives an overview of the model the analyst is constructing. (B) Based on the variable relationships analysts specify, Tisane infers candidate main effects that may be potential confounders. Tisane asks analysts if they would like to include these variables, explaining in a tooltip (C) why the variable may be important to include. (D) Tisane only suggests interaction effects if analysts specify moderating relationships in their specification. This way, Tisane ensures that model structures are conceptually justifiable. (E) From the data measurement relationships analysts provide, Tisane automatically infers and includes random effects to increase generalizability and external validity of statistical findings. (F) Tisane assists analysts in choosing an initial family and link function by asking them a series of questions about their dependent (e.g., Is the variable continuous or about count data?). To help analysts answer these questions and verify their assumptions about the data, Tisane shows a histogram of the dependent variable.

because analysts may not have specific hypotheses involving interactions (*DG3 - Guidance and control*). If analysts do not specify any moderating relationships, Tisane does not suggest any interaction effects, preventing analysts from including arbitrary interactions that may be conceptually unfounded (*DG1 - Conceptual knowledge, DG2 - Validity*).

In the data clustering tab, Tisane shows analysts which random effects it automatically includes based on the selected main and interaction effects. Unlike main and interaction effects, Tisane automatically includes random effects in order to maximize model generalizability (*DG2 - Validity*). If there is a random slope and random intercept pertaining to the same unit, Tisane asks analysts if they should be correlated or uncorrelated. We provide this option because analysts may have relevant domain expertise to make this decision (*DG3 - Guidance and control*). By default, Tisane correlates the random slope and random intercept.

The final tab, data distribution, helps analysts examine their data and select an initial family and link function to try. Appropriate selection of family and link functions depends on the data type of the dependent variable and the distribution of model residuals. Therefore, the selection can only be assessed after choosing a family and link function in the first place.

For an initial statistical model to consider, Tisane narrows the set of family functions considered based on the declared data type of variables (see 5.3.2) and lightweight viability checks, such as ensuring that a Poisson distribution is only applicable for variables that have nonnegative integer values. Tisane asks questions designed to uncover more semantically meaningful data types (e.g., counts) than are provided at variable declaration. Analysts without data can answer these questions as they are planning their studies (*DG4 - Statistical planning*). For the selected family candidate, Tisane automatically selects the default link function based on the defaults for `statsmodels` [PSTsd20] and `pymer4` [Jol18]. Analysts can then choose a different link function, as long as it is supported.

Output

There are two outputs of the interactive compilation: (ii) an executable modeling script and (ii) a log of GUI choices. To increase transparency of the authoring process, Tisane provides a log of user selections in the GUI as documentation, which the analyst can include in pre-registrations, for example (*DG4 - Statistical planning*). In the output script, Tisane includes code to fit the model and plot residuals against fitted values in order to assess the appropriateness of family and link functions, as is typical when examining family and link functions. The output script also includes a comment explaining what to look for in the plots and an online resource for further reading. Should analysts revise their choice of family and link functions, they can re-generate a script through the Tisane GUI.

5.4 Initial evaluation: Case studies with researchers

Given Tisane’s novel focus on deriving and guiding analysts toward valid statistical models, we assessed how Tisane affects data analysis practices in three case studies with researchers. The following research questions guided the evaluation:

- **RQ1 - Workflow** How does Tisane’s programming and interaction model affect how analysts author models? Specifically, what does Tisane make noticeably easier or more difficult when conducting an analysis?
- **RQ2 - Cognitive fixation** Where do researchers report spending more time or attention when using Tisane? How does this compare to their fixation during analyses typically?
- **RQ3 - Future possibilities** When do researchers imagine using Tisane in future projects, if at all? What additional support do researchers want from Tisane?

We recruited researchers through internal message boards and individual contacts. We intentionally recruited researchers at different stages of the research process—study planning, data analysis for publication, and ongoing model building and maintenance. We believed this could help us more holistically evaluate Tisane’s impact on data analysis. We met with researchers over Zoom (R1, R3) and in person (R2) to discuss their use cases, observe them use Tisane for the first time, and ask for open-ended feedback. We pointed researchers to the Tisane tutorial for installation instructions and examples but otherwise encouraged the researchers to work independently. We answered any questions researchers had while using Tisane. Each study session lasted approximately 2 hours. At the end, two of the three researchers (R1, R3) said they planned to use Tisane again over the next two months.

5.4.1 Case Study 1: Planning a new study

R1, a clinical psychology PhD student, had recently submitted a paper and was planning a follow-up. R1 reported that she had never taken a formal class on modeling techniques but taught herself for her last paper. Her general workflow involved consulting with and mirroring what others in her research group did even if she did not completely understand why. R1 did not program often but said she had “enough coding experience to understand this kind of...[sample program].” Although familiar with Python, R1 preferred M+ [Mut23] and SPSS [SPS21]. She was interested in using Tisane to brainstorm new studies and research questions.

Using Tisane. After installation, R1 read through one of the computational notebook examples available in the Tisane GitHub repository. While reading, R1 asked clarifying questions about the variable types and syntax. R1 explained that the Design class felt novel because she had

never seen the concept of a study design in data analysis code before. When the first two authors explained that it was supposed to be the equivalent of the statement of a study design in a paper, R1 remarked that usually, she “[kept] that in [her] head, which [she] probably shouldn’t” (**RQ2 - Cognitive fixation**). Without a concrete data set, R1 preferred to walk through more examples rather than author a script of her own.

While reading an example, R1 drew a parallel between the tabs in SPSS dialogs for specifying models and the tabs in the Tisane GUI, noting that SPSS had a tab for control variables. R1 also wanted the ability to distinguish between “control variables” and other independent variables in the Tisane GUI. R1 explained that this would map more closely to how psychologists think about analyses. Future work could incorporate additional language constructs, such as a new data type for controls, for different groups of users (**RQ3 - Future possibilities**).

At the end of the study session, R1 remarked how Tisane “fills in a lot of the...gaps” in data analysis (**RQ1 - Workflow, RQ2 - Cognitive fixation**). The first gap R1 discussed was the *programming gap* between scientists and statistical tools. R1 believed that, for scientists who were not comfortable with programming, “they should probably be running less complex models, or first learn how to code” even if the complex models would be most appropriate. The second gap R1 discussed was the *statistical knowledge gap* in tools. R1 explained that in her experience, R provides support for more complex models but little guidance for what those models or statistical tests should be, requiring “top down assumption[s].” Thus, to R1, Tisane bridged the gap between tools like SPSS and R by requiring minimal programming and providing modeling support. Put another way, Tisane bridged the gulf of execution [Nor13] for R1 that previous tools had not.

5.4.2 Case Study 2: Analyzing data for a paper submission

R2, a computer science PhD student, had conducted a within-subjects study where 47 participants used four versions of an app for one week each (four weeks total). The motivating research question was how the different app designs led to psychological dissociation. Although R2 had expected to collect multiple survey responses for each participant each day, they only had aggregate daily self-report measures due to an error in the database management system. In the past, R2 reported having extensively explored their data and consulting others, but for this paper, they had not explored their data prior to fitting models because they felt more confident in their modeling skills. For analyses, R2 preferred R but had general Python programming experience. Prior to using Tisane, R2 had authored linear mixed effects models in R for their study. They were interested in using Tisane to check their analyses prior to submitting their paper to CHI.

Using Tisane. R2 wrote their scripts by adapting an example from the Tisane GitHub repository. As R2 considered which conceptual relationships to add, they reasoned aloud about if they

should state causal or associative relationships between various measures and dissociation (**RQ2 - Cognitive fixation**). After some deliberation, they said, “I don’t feel comfortable [making a causal statement],” and instead specified `associates_with` relationships. R1’s hesitation to assert causal relationships confirms prior findings that specifying formal causal graphs is difficult for domain researchers [SSY20; SV18; VDN⁺13] and our design choice to allow for association edges. In addition, R2 was initially unsure about how to specify the `number_of_instances` for their measures since their original study design was unbalanced. After asking for clarification about `number_of_instances`, R2 declared all the measures with the parameter `number_of_instances` set equal to date.

Next, R2 ran their script and used the Tisane GUI in a browser window. Based on Tisane’s recommended family and link functions, R2 realized the models they had previously authored in R using a Gaussian family were inappropriate. Due to a bug that we have since fixed, Tisane suggested a Poisson family that R2 used to generate a script, but this was an invalid choice given that not all dependent variable values were nonnegative integers. R2 explored other family distributions and generated a new script using an Inverse Gaussian family. When executed, the second output script issued an error due to the model inference algorithm failing to converge. R2 made a note to look into this model further on their own.

Once finished using Tisane, R2 commented that their analysis with Tisane was more streamlined (**RQ1 - Workflow**) in contrast to their very first paper where they had tried “every single kind of model that [they] could” until finding “the one that fits best,” even if it was “one that no one would have heard of.” R2 also stated they would be interested in using Tisane earlier in their analysis process in the future (**RQ3 - Future possibilities**). Based on their experience with Tisane, R2 questioned their previously authored linear mixed effects model, and said it was “unnerving” to discover an issue so close to a deadline. At the same time, they expressed, “if it’s incorrect, I should know before I submit.” A day after the study, R2 contacted the authors to inform them that they had decided to update their analyses from linear mixed effects models to generalized linear mixed effects models. They reported using the Inverse Gaussian family after visualizing and checking the distribution of residuals with help from the output Tisane script. The Inverse Gaussian family was appropriate because their dependent variable’s values were all nonnegative and displayed a slight positive skew. R2’s experience with Tisane suggests that Tisane can help researchers catch errors and lead them to re-examine their data, assumptions, and conclusions.

5.4.3 Case Study 3: Developing models to inform future models

Employed on a research team, R3 analyzes health data at the county, state, and national levels to estimate health expenditure and inform public policy. R3 develops initial models that are used

to validate and generate estimates for larger, more comprehensive models. Due to the scale of data and established collaborative workflows, R3 typically works in a terminal or RStudio through a computing cluster and had very little experience with Python. Despite working on statistical models every day, R3 described himself as “not...a great modeler.” R3 was interested in using Tisane to determine what variables to include as random effects in a model.

Using Tisane. R3 used Tisane in a local Jupyter notebook as well as on his team’s cluster. R3 used the Tisane API overview reference material on GitHub to start writing his program, which involved copying and pasting the functions with their type signatures and then modifying them to match his dataset and incrementally running the program. The most common mistake R3 made while authoring his Tisane program was to refer to variables using the string names in the dataset (e.g., "year") instead of the variable’s alias (e.g., year_id), an idiom common in R but not in Python.

While authoring his Tisane program, R3 found the number_of_instances parameter redundant, especially because his data is always “square.” Every state_name in his data set had 30 rows of data, corresponding to the year_ids 1990-2019. This is in contrast to R2, whose study design was unbalanced and resulted in variable numbers of observations per participant that needed to be aggregated. Based on R3’s feedback, we added functionality to infer number_of_instances for each unit, which analysts can inspect by printing the variable.

While giving open-ended feedback on Tisane, R3, similar to R1, liked how Tisane helped “fill [the] gap in...[his] knowledge” (**RQ2 - Cognitive fixation**). Given the diversity of models R3 works with, R3 found Tisane’s focus on GLMs and GLMMs a “little limiting” and also wished to make Tisane “run without...the mouse” in a script, as is typical in his workflow (**RQ1 - Workflow**). Specifically, R3 described how he and his collaborators typically want to explore a space of models and run them in parallel. Nevertheless, R3 foresaw using Tisane in three types of modeling tasks common in his work: (i) exploratory modeling to determine if there are any interesting relationships between variables, (ii) authoring and comparing multiple models for prediction, and (iii) working out the precise model specification after identifying variables of interest (**RQ3 - Future possibilities**).

System changes and Takeaways

We fixed bugs and iterated on Tisane’s GUI based on feedback from researchers. The largest change we made was to the data distributions tab. The data distributions tab we tested with researchers visualized the dependent variables against simulated distributions of family functions and included the results of the Shapiro-Wilk and D’Agostino and Pearson’s normality tests. All three researchers reported becoming more aware of their data due to the visualizations. However, researchers’ enthusiasm for the feature made us wary that visualizing the simulated data could

mislead less careful analysts to believe that family and link functions pertain to variable distributions rather than the distributions of the model’s residuals. To avoid such errors while still helping analysts become more aware of their data, we removed the simulated visualizations and normality tests and instead provide questions about the semantic nature of the dependent variable collected, as discussed in Section 5.3.2.

Overall, Tisane streamlines the analysis process (**RQ1 - Workflow**) in part because researchers report formalizing their conceptual knowledge into statistical models more directly (R1, R2). Although Tisane does not eliminate the need for model revision, Tisane may scope the revisions analysts consider to significant issues instead of details that may detract from the analysis goals (R2). Additionally, researchers reported a perceived shift in their attention from keeping track of and analyzing all possible modeling paths to their research questions and data assumptions (**RQ2 - Cognitive fixation**) while planning a new study and analysis (R1) as well as while preparing a research manuscript (R2). Future adoption of Tisane may depend on the complexity of analyses (**RQ3 - Future possibilities**) (R3). For instance, Tisane may provide a streamlined alternative to false starts due to misspecifications for simpler analyses (R1, R2, R3). For more complex models and studies, Tisane may act more as a prototyping tool for statistical models, helping researchers start at a reasonable model that they can then revise (R2, R3).

5.5 Limitations and Motivation for Re-design

While overall positive, the case study made us aware of confusing keywords and language constructs in Tisane. In order to improve Tisane and probe more closely into what challenges statistical novices face when expressing their domain knowledge, we engaged in an iterative process to re-design Tisane. We started with a lab study using Tisane [JSHJ22] to elicit statistical non-experts’ implicit definitions and assumptions about Tisane’s keywords and identify opportunities to refine Tisane’s DSL and interactivity.

5.6 Elicitation lab study

Our aim was to increase the expressivity of rTisane to represent analysts’ implicit domain knowledge. We used the first release of Tisane [JSHJ22] to probe analysts’ internal processes and derive design goals Subsection 5.6.3 for re-designing Tisane.

5.6.1 Method

We recruited participants through a graduate-level quantitative research methods course as a convenience sample to control recent exposure to statistical concepts. Five computer science PhD students participated.

The study consisted of two parts: (i) a take-home assignment and (ii) an in-lab session. The take-home assignment asked participants to read a recently published CHI paper [WWL21]⁶ and describe the paper’s research questions and hypotheses, the authors’ conceptual models, the study’s design, and ways to analyze the data to answer the research questions. The assignment was designed to ensure that participants engaged with the paper’s key ideas before coming into the lab. The researcher reviewed each submission to prepare participant-specific questions for a semi-structured, think-a-loud lab session.

At the start of the lab session, participants reviewed their homework submission to remind themselves of the paper. The paper and participants’ homework responses remained available for reference throughout the study. Then, participants completed three tasks: (i) declaring variables, (ii) expressing conceptual models, and (iii) specifying study designs. For each task, participants started with Tisane’s language constructs to express their intent and discussed their confusions, how they understood each presented construct, and what they wanted to specify but could not (if applicable).⁷ The researcher repeatedly reminded participants that the constructs presented were prototype possibilities and that expressing their intentions was more important than using the constructs or getting the syntax correct. Throughout, the researcher paid particular attention to where Tisane broke down for participants and asked follow-up questions to probe deeper into why. The researcher considered such breakdowns as openings into semantic mismatches between the end-user and the DSL.

We iteratively coded homework submissions, audio transcripts from the lab study sessions, and lab study artifacts. We also consulted the researcher’s detailed notes from the lab sessions.

5.6.2 Key Observations

All participants demonstrated a working knowledge of the assigned paper’s motivating research questions, study design, and general study procedure. We made four key observations about what and how statistical non-experts wanted to express their conceptual models: using varying degrees of specificity, separating moderation from bivariate relationships, distinguishing between known

⁶We chose the specific paper because we believed its topic (i.e., empathetic biosignals) would be broadly approachable and the statistical methods used (i.e., generalized linear models) would be familiar with students enrolled in the research methods course.

⁷For the lab study, we re-implemented Tisane (originally in Python) in R since (i) R is a widely used programming language for data science and (ii) the research methods course taught and used R.

and hypothesized relationships, and considering alternative conceptual models. Participants also suggested syntactic sugar options to improve the DSL's usability. Based on these observations, we derived design goals for re-designing Tisane (see Section 5.7).

Participants express conceptual knowledge with varying details.

Contrary to the popular belief that higher levels of abstraction are better for end-users, we found that statistical non-experts want to move up and down the ladder of abstraction when expressing conceptual models.

When defining “causes,” P2 described “[Causes] is...like when we teach logic...it’s like implication, right?....So I’m saying if we are observing an emotion and...emotion observed can lead to a change in emotional perspective.” P0, P1, and P3 contrasted a bidirectional relationship between variables, formerly encapsulated in the `associates_with` construct in Tisane, to their implicit understanding of “causes.” For instance, P1 stated “the most like, utilitarian definition by if A causes B, then by changing A, I can change B whereas `associates_with` means that...if I can turn dial A, B might not change.” In addition to differentiating between causal and associative relationships, three participants [P0, P1, P3] provided statements of *specifically how* a variable influenced another in the conceptual models submitted as homework. For example, P0 wrote, “Hearing a heartbeat that seems to be aligned with visual cues makes someone feel *more* strongly what another person is feeling” (emphasis added), specifying a positive influence of “hearting a heartbeat” on empathy. These observations suggest that analysts have an intuitive understanding of causality but bluntly stating that a variable causes another does not capture the richness or nuance of their implicit domain knowledge. Additional annotations about how a variable influences another are necessary.

Participants find moderation difficult to separate from bivariate relationships.

Participants consistently found Tisane’s `moderates` construct difficult to understand [P0, P1, P2, P3]. Participants expressed confusion about what moderation implied about the relationship between two variables. For example, P3 grappled with if “moderates” was shorthand for expressing associative relationships between each independent variable and the dependent variable, how moderation implies causal relationships, and if statistical and conceptual definitions of moderation differed from each other: “[L]et’s say there’s two independent variables and one dependent variable. And each of the [independent] variables individually is not correlated with the outcome. But if you put them together, then the correlation appears....I mean, it’s sort of a philosophical question of whether, like each of the ones individually causes [the dependent variable] in that case. But thinking from a...statistical perspective, I think that’s a situation where you might be able to express...language and experience level together cause lines of code but individually they don’t because no individual

correlation would appear there.” Therefore, a clear delineation between bivariate relationships and partial statistical specifications of interaction terms is necessary.

Participants distinguish between known and suspected relationships.

Participants described relationships established in prior work as “assumptions” or “assertions” to check separately from the key research questions that tested “suspected” relationships. P0 described how “maybe we have to differentiate as to like the *known* [relationships] are kind of the things you’re *assuming* there’s relationships between these things whereas the *suspected...[are]* the things kind of like your research questions are saying like, ‘We *think* there’s this relationship but...it’s what we’re testing for’ (emphasis added). Similarly, P4 suggested that Tisane should warn end-users when assumptions about known relationships are violated in a given data set: “I would also say that it would be very handy to be able to say, kind of *assert* that language has no effect on the line of code. And be warned if it’s not the case, like if your *assertion* is not...verified automatically with the DSL, but warned...that while your *assumption* is not holding there is actually an effect, which could be very handy on your study.” (emphasis added). The inability to indicate relationships that are either known or suspected in Tisane may explain why analysts repeatedly preferred less technical verbs, such as “influences” [P0] or “leads to” [P3]. For instance, P0 explained how she preferred “influences” over “causes” because “I guess it’s like *a level of sureness* in it in which, like, ‘cause’ feels more confident in your answers than ‘influences’” (emphasis added). Providing a way to label conceptual relationships as assumptions or the focus of the present analysis could make `causes` and `associates_with`, the bivariate relationships in Tisane, more approachable.

Participants consider alternative conceptual structures in the face of ambiguity.

Participants grappled with what specific structures in a conceptual model meant. P1 and P3 described how a bidirectional relationship between two variables were really due to hidden, confounding variables causing both variables. P3 described how “in the real world, when, when these bidirectional things happen, it means there’s sort of this middleman complex system. Or some like underlying process of which [two variables are] both components...” Another participant, P2, wondered aloud about how even what appears like a direct relationship, may actually be a chain of indirect or mediated relationships at a lower granularity: “It’s like Google Maps. If you zoom out enough, that arrow becomes a direct arrow.” These observations suggest that while participants can deeply reflect on what could be happening between variables conceptually, they need help exploring and figuring out which of these structures matches their implicit understanding.

Participants expected more syntactic sugar for specifying data collection details.

While our focus was on improving the support for conceptual modeling, we made a few observations about challenges analysts faced when specifying data collection details. First, analysts expected *experimental conditions to be standalone concepts*. In Tisane, experimental conditions can be specified as a Measure of a Unit. Instead, P0 and P4 had separate conceptual categories for conditions and measures in their mental models of study designs. P4 preferred a separate condition data type currently unavailable in Tisane because the term “Measure” did not create a “bucket” appropriate for conditions. Second, participants were interested in specifying trials, stimuli, and responses elicited during each trial alongside participants: “I want to have a trial unit that is nested within trials, which is nested within or maybe I could just have trial nested within Participant, but I’m not seeing a way to clearly delineate or like to denote that” [P1]. Future work should more closely examine and iterate on language constructs and idioms for representing data collection procedures.

5.6.3 DSL Re-design Goals

Based on our lab study observations, we derived four design goals for re-designing Tisane’s DSL to more accurately capture analysts’ implicit conceptual models:

- **DG1 - Optional specificity:** Analysts should be able to provide optional details about how variables change in relation to each other (e.g., positive or negative changes in values) when describing conceptual relationships.
- **DG2 - Interactions as partial specifications:** Analysts should annotate conceptual models with interaction terms they want to include in an output statistical model.
- **DG3 - Consideration of possibilities:** When expressing ambiguous relationships, analysts should have support in considering and picking among multiple possible conceptual structures.
- **DG4 - Distinction between assumed and hypothesized:** Analysts should be able to distinguish between assumed and hypothesized relationships in their conceptual models.

In the second release of Tisane, we addressed these goals through new language constructs. We also supported syntactic sugar to more accurately capture study design details.

5.7 Second Release: rTisane

rTisane consists of (i) a DSL for analysts to express their conceptual models and (ii) interactive disambiguation steps to compile this high-level specification into a script fitting a statistical model.

So far, we have implemented rTisane for GLMs. Given the breadth of findings from the elicitation lab study, we narrowed the scope from Tisane in order to really focus on designing and testing a set of language constructs core to conceptual modeling. There are two key challenges to designing rTisane: (i) ensuring the DSL's constructs can express analysts' implicit conceptual models accurately and (ii) balancing usability with rigor, allowing analysts to express their often “fuzzy” conceptual assumptions without losing precision to derive a statistical model.

```

1 library(rTisane)
2
3 # Declare variables
4 # Person: Observational unit
5 person <- Unit(name = "person")
6 # Age: Continuous measure
7 age <- continuous(unit = person, "Age")
8 # Race, 5 categories:
9 # White, Black/African American, American Indian or Alaska Native, Asian or Pacific Islander, Mixed Race
10 race <- categories(unit = person, "Race", cardinality = 9)
11 # Highest Education Completed, 5 ordered categories
12 edu <- categories(unit = person, "Education", order=list("Grade 12", "1 year of college", "2 years of
    college ", "4 years of college ", "5+ years of college "))
13 # Current Employment Status, 3 categories: N/A, Works for wage, Self-employed
14 employ <- categories(unit=person, "Employment", cardinality=2)
15 # Sex, 2 categories: Male, Female
16 sex <- categories(unit = person, "Sex",cardinality = 2)
17 # Income: Continuous measure
18 income <- continuous(unit = person, "Income")
19
20 # Construct a conceptual model
21 cm <- ConceptualModel() %>%
22   assume(causes(age, income)) %>%
23   assume(causes(race, income)) %>%
24   hypothesize(relates(edu, income)) %>%
25   hypothesize(relates(age, edu)) %>%
26   hypothesize(relates(race, edu)) %>%
27   hypothesize(relates(sex, edu)) %>%
28   hypothesize(relates(employ, income)) %>%
29   hypothesize(causes(sex, income)) %>%
30   interacts (race, sex, dv = income) %>%
31   interacts (age, edu, dv = income)
32
33 # Query for a statistical model

```

```
34 query(conceptualModel=cm, iv=edu, dv=income)
```

Listing 5.1: Sample rTisane program adapted from P8's in the summative evaluation. Specifying cardinality is optional with data.

5.7.1 rTisane's DSL

Like Tisane, analysts express variables, a conceptual model, and a query for a statistical model. rTisane's DSL prioritizes expressivity and usability

Declaring variables

Analysts can express two types of variables: Units and Measures. Units represent observational or experimental units from which analysts collect data (see Line 5 in Listing 5.1). A common unit is a participant in a study, so rTisane provides syntactic sugar for constructing a Participant unit directly. Participant is implemented as a wrapper for declaring a Unit.

Measures are attributes of Units collected in a data set, so they are declared through a Unit. Measures can be one of four types: continuous, unordered categories (i.e., nominal), ordered categories (i.e., ordinal), and counts (see Lines 7-18 in Listing 5.1). rTisane provides syntactic sugar for declaring Conditions as either unordered or ordered categories. Analysts declare unordered and ordered categories through the categories function. Analysts can specify a variable is ordered by passing a list to the order parameter. Otherwise, the variable is considered unordered. Analysts can use continuous and count functions to declare continuous and count Measures.

Expressing a conceptual model

Once analysts have constructed variables, they can specify how these variables relate conceptually. To do so, they construct a ConceptualModel and add variable relationships to it (Lines 21-31 in Listing 5.1). The conceptual model is represented as a graph where the variables are nodes and the relationships are edges.

There are two types of relationships: causes and relates. causes indicates a unidirectional influence from a cause to an effect. causes introduces a directed edge from the cause node to the effect node. relates indicates that two variables are related but exactly how is ambiguous because the analyst is uncertain about the direction of influence. relates introduces a bi-directional edge between two variables. During a disambiguation step, rTisane will walk analysts through possible graphical structures that a bi-directional edge could represent (**DG3 - Consideration of possibilities**). To derive a statistical model, rTisane requires an analyst to assume a direction of influence.

Towards the design goal of **DG1 - Optional specificity**, rTisane allows analysts to optionally specify when and then parameters in the causes and relates functions. There are four comparisons analysts can specify in when and then: increases (for continuous, ordered categories, counts), decreases (for continuous, ordered categories, counts), equals (for any measure type), and notEquals (for any measure type). Supporting optional specificity is designed to (i) make the rTisane program an accurate document of analysts' implicit assumptions and (ii) suggest ways to resolve conceptual ambiguity during disambiguation (**DG3 - Consideration of possibilities**).

To add relationships to the conceptual model, analysts must assume or hypothesize a relationship. This distinction supports how analysts distinguish between assumed, or strongly held, and hypothesized, or more uncertain, relationships. rTisane requires analysts to make these explicit distinctions (**DG4 - Distinction between assumed and hypothesized**) when adding conceptual relationships to a conceptual model. In addition to specifying a relationship type, analysts must either assume or hypothesize a relationship.

Analysts can also specify interactions between two or more variables by declaring interacts. Interactions are annotations to conceptual models and are added to the graph without assume or hypothesize. Interactions provide additional information about existing relationships in the conceptual model (**DG2 - Interactions as partial specifications**).

Querying for a statistical model

Analysts query rTisane for a statistical model based on the input conceptual model (Line 31 in Listing 5.1). The query asks for a statistical model to accurately estimate the average causal effect (ACE) of the independent variable on the dependent variable. The querying process initiates the interactive compilation process and results in an R script specifying and fitting a generalized linear model. During interactive compilation, analysts engage in two loops to disambiguate their (i) conceptual model and (ii) output statistical model.

5.7.2 Two-step interactive compilation

There are two phases to interactively compiling a conceptual model to a statistical model: (i) conceptual model disambiguation and (ii) statistical model disambiguation. We added conceptual model disambiguation to address the need to explore possible conceptual structures for resolving ambiguities introduced by relates (**DG3 - Consideration of possibilities**).

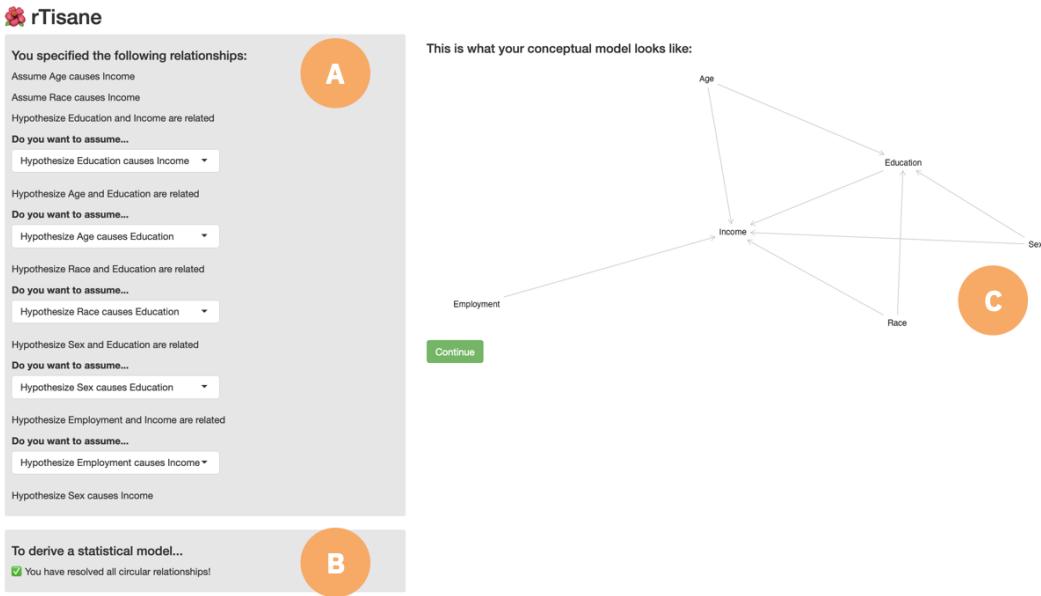


Figure 5.3: rTisane’s conceptual model disambiguation interface. A. Options for resolving ambiguities in the conceptual model due to *relates* relationships. B. Check and follow-up questions for breaking any cycles that hinder statistical model derivation. C. Graph visualizing expressed conceptual model.

Conceptual Model Disambiguation

The goal of conceptual model disambiguation is to make analysts’ expressed conceptual models precise enough to derive a statistical model, achieving usability and rigor. Conceptual model disambiguation involves breaking cycles in the conceptual model by (i) picking a direction for any *relates* relationships and/or (ii) removing edges. Cycles are necessary to break because they imply multiple different data generating processes that could lead to different statistical models. In this way, conceptual model disambiguation can help analysts reflect on and clarify their implicit assumptions.

To disambiguate conceptual models, rTisane uses a GUI. Figure 5.3) shows the conceptual model disambiguation interface for the input program in Listing 5.1. The GUI shows a graph representing analysts’ conceptual models. If there are any *relates* relationships, rTisane suggests ways analysts could assume a direction of influence. Additionally, rTisane suggests ways to break any cycles in the conceptual model. As analysts make changes, the visible graph updates. The GUI also explains why both these steps are necessary to derive a statistical model.

Once analysts have disambiguated their conceptual models, rTisane updates the internal graph representation and derives a space of possible statistical models. To narrow this space of possible statistical models down to one output statistical model, rTisane asks additional follow-up disam-

biguating questions.

Statistical model derivation and disambiguation

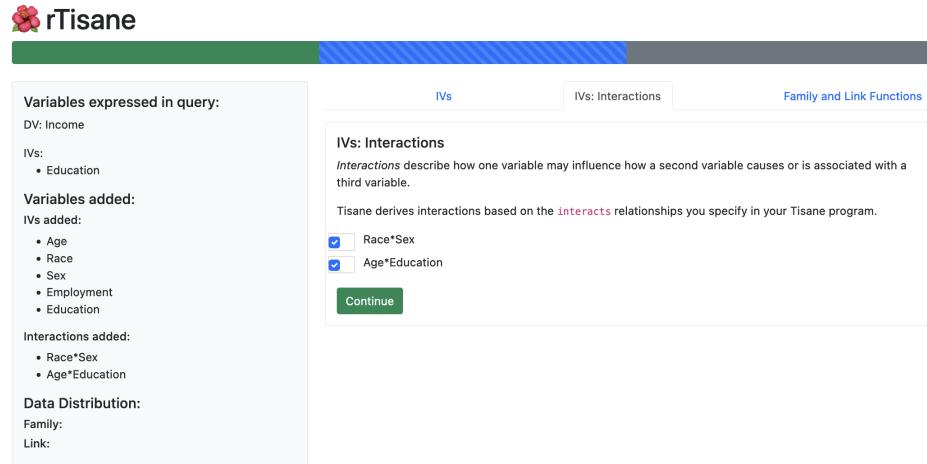


Figure 5.4: rTisane’s statistical model disambiguation interface.

To formulate possible statistical models, rTisane considers potential covariates to control for confounding, interactions, and family and link functions.

To determine confounders, rTisane uses more recent recommendations from Cinelli, Forney, and Pearl [CFP20]⁸. Cinelli et al.’s recommendations are based on a meta-analysis of studies examining the impact of confounder selection based on graphical structures on statistical modeling accuracy. By following Cinelli et al.’s recommendations, rTisane includes confounders that help assess the average causal effect of the query’s independent variable on the dependent variable as accurately as possible.

rTisane searches for interactions analysts annotated in their conceptual models and suggests any involving the query’s dependent variable. Otherwise, rTisane does not consider any interactions.

rTisane determines family and link functions based on the query’s dependent variable data type. Because rTisane compiles down to statistical models fit using the `lme4` package in R, rTisane is limited to the family and link functions supported in `lme4`. For instance, for queries involving continuous dependent variables, rTisane considers Gaussian, Inverse Gaussian, and Gamma families. For counts, rTisane considers Poisson and Negative Binomial families. For ordered categories, rTisane considers Binomial, Multinomial, Gaussian, Inverse Gaussian, and Gamma family functions. For unordered categories, rTisane considers Binomial and Multinomial family functions. rTisane considers any link functions `lme4` supports for these family functions.

⁸Tisane relied on Vanderweele’s recommendations for confounder selection [Van19], but in rTisane we opted for more recent recommendations

In the GUI, analysts have the option to remove any confounders or interactions based on their domain knowledge. Based on prior experience or domain recommendations, analysts can also pick a family and link function pair if multiple possibilities could apply.

5.8 Summative Evaluation: Controlled lab study

Two research questions motivated our evaluation of rTisane:

- **RQ1 - Conceptual models** What is the impact of rTisane on conceptual modeling?
- **RQ2 - Statistical models** How does rTisane impact the statistical models analysts implement? Specifically, how well do the statistical models analysts author on their own vs. with rTisane fit the data? How are their formulations similar or different?

5.8.1 Study design

We conducted a within-subjects (Tool support: rTisane vs. none) think-a-loud lab study that consisted of four phases. We designed the study based on the assumption that conceptual modeling is a helpful strategy when specifying statistical models. As a result, all participants completed the phases in the following order.

- **Phase 1: Warm up.** We presented participants with the following open-ended research question: “What aspects of an adult’s background and demographics are associated with income?” We asked participants to specify a conceptual model including variables they thought influenced income. This warm-up exercise helped to externalize and keep track of participants’ pre-conceived notions and assumptions prior to seeing a more restricted data schema.
- **Phase 2: Express conceptual models** We presented participants with a data schema describing a dataset from the U.S. Census Bureau. We then asked participants to specify a conceptual model using only the available variables. At the end, we asked participants about their experiences specifying their conceptual models in a brief survey and semi-structured interview.
- **Phase 3: Implement statistical models** We asked participants to implement “a statistical model that assesses the influence of variables [they] believe to be important (in the context of additional potentially influential factors) on income,” relying on only their conceptual model. We then asked participants about their experiences implementing statistical models through a brief survey and semi-structured interview.

- **Phase 4: Exit interview.** The study concluded with a survey and semi-structured interview where we asked participants about their experience in the study, reactions to using rTisane, and connecting conceptual models to statistical models.

In order to assess the effect of tooling on conceptual models and the quality of statistical models, we counterbalanced the order of tool support, or if participants completed each task with or without rTisane first. The order of tool use was the same for Phase 2 and Phase 3. Within each of Phase 2 and Phase, half the participants completed the task on their own (without rTisane) then with rTisane. The other half started with rTisane and then did the task on their own. Prior to using rTisane in Phases 2 and 3, participants followed a tutorial introducing the relevant language constructs. Section C.1 contains all the study materials.

Participants We recruited 13 data analysts on Upwork. We screened for participants who reported having experience with authoring generalized linear models and using R at a three or higher on a five-point scale. Table 5.3 summarizes the participants' backgrounds.

Update based on recruitment: On average, participants reported having conducted N projects using generalized linear models.

All studies were conducted over Zoom. Participants used rTisane on a remote controlled computer, so they did not have to install it on their own. Each study lasted between two and three hours. Participant was compensated \$25 per hour. We recorded participants' screens, video, and audio throughout the study. We then transcribed the audio and used detailed researcher notes for qualitative analyses.

| Participant | Field | Job role | Reported experience | Familiarity with GLM | Comfort using R |
|-------------|-------|----------|---------------------|----------------------|-----------------|
| P1 | | | | | |

Table 5.3: Participants in summative evaluation.

5.8.2 Analysis Approach

Our analysis procedure consisted of two parts: (i) a thematic analysis of lab notes, transcripts, and open-ended survey questions and (ii) an artifact analysis of conceptual models and statistical models analysts authored with and without rTisane. For the conceptual models, we compared their form and content between tool support conditions. For the statistical models, we compared the overall statistical approach, specific statistical model formulation, and rationale for analysis decisions and conclusions. We also compared two goodness of fit measures between statistical models: AIC and BIC. We iterated on the thematic analysis and artifact analysis separately at first and then interpreted emergent observations across the two analyses.

One of the 13 participants dropped out part way through the study due to discomfort with programming in front of the researchers. We analyzed the data we were able to collect from them.

5.8.3 Findings

RQ1: rTisane's Impact on Conceptual Models

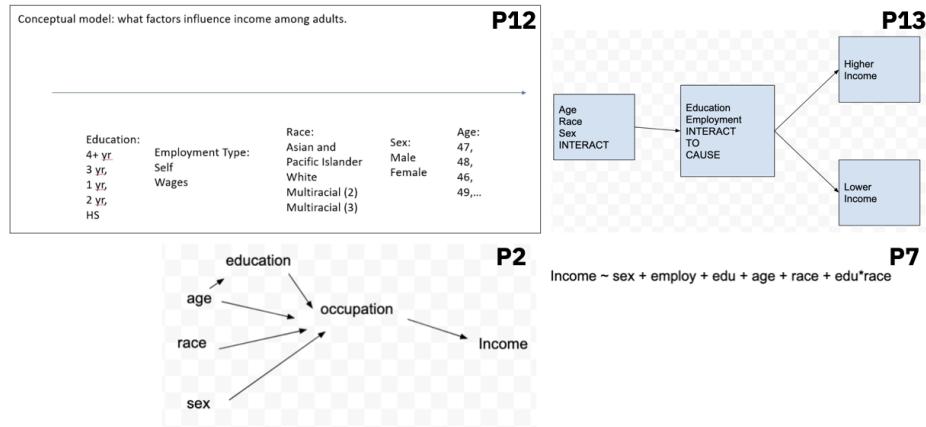


Figure 5.5: Example conceptual models from participants in the summative evaluation without using rTisane.

Key takeaway: rTisane scaffolded and productively constrained how analysts expressed their conceptual models. As a result, analysts reflected on implicit domain assumptions more deeply, considered new relationships, and felt they accurately externalized their implicit assumptions.

The conceptual models analysts expressed on their own were diverse in form, meaning/content, and complexity. The majority [P2, P4, P5, P8, P11, P13] invoked a graph-like structure. [P2, P4, P8 used rTisane second; P5, P11, P13 used rTisane first]. Figure 5.5 illustrates four example conceptual models from participants.⁹ Participants also described their conceptual models verbally [P10], in natural language text [P6, P9], and as a timeline [P12]. P7, who used rTisane first, even jumped to expressing their conceptual model in a statistical model. P12's conceptual model was particularly creative. His timeline featured variables ordered starting on the left by how much an individual could intervene upon them (Figure 5.5). P12's conceptual model reiterates our finding from the exploratory lab study that analysts want to capture nuanced meaning in a conceptual model.

⁹An example conceptual model given in the task instructions may have biased analysts towards a graphical structure.

Ten participants involved all five independent variables from the data set in their conceptual models [P2, P3, P4, P5, P7, P8, P9, P11, P12, P13]. Two participants [P7, P13] also included interactions between variables in their conceptual models. For instance, P13 specified a complex conceptual model (Figure 5.5) where age, race, and sex interacted to cause an interaction between education and employment, which then causes income.

Theme: Without rTisane, analysts found it difficult to express conceptual nuances.

In a survey and interview about their conceptual modeling experiences, participants shared that they found it difficult to author conceptual models without tool support due to doubts about how to communicate nuances in relationships [P3, 13] and concerns about mis-specifying relationships beyond their domain knowledge [P5, P10]. P13 explained how they wanted to “[i]dentify how I may weigh certain variables based on my general awareness and knowledge and overall weights of each variable of how one may affect income more or less in various circumstances.” Similarly, P8 described specifying their conceptual model as a general “struggle” because “When doing it myself, there are so many possibilities [of expression].” While rTisane is not designed to prevent mis-specifications due to limited domain knowledge, we found that rTisane’s formalism removed the need for analysts to come up with how to express their domain knowledge. They could focus on expressing what they knew.

Theme: rTisane encouraged analysts to think about their domains more deeply.

rTisane ’s DSL deepened participants’ thinking [P3, P4, P7, P8, P10, P12, P13], giving them, as P12 described, a structure to explore the “boundaries of their domain knowledge.” P3 explained how even after specifying conceptual models on her own, rTisane ’s four composable relationships (`assume/hypothesize x causes, relates`) facilitated a deeper consideration of each relationship and what she knew about each: “Having to think about specifics like ’Do we know the direction of the relationship’ or ’What happens when a category increases/decreases’ actually helped me put my thoughts out more clearly. I was able to think about more possible scenarios that could conflict with my current assumption, which I was probably not doing [before]...In conclusion, I want to say that looking at four possible ways to write a relationship made me think more about each one of them.” [P3] Similar to P3, P10 explained, “My thinking was that before I didn’t have much idea about how can I link my variable with the output [variable], and how this can interact. And so it may need some trial and error... using this API, there are predefined functions, they are translated in R language, cause or relates, it made my task easier. This translation was not on me anymore.” Furthermore, P4 explained how the DSL’s support for optional specificity “encouraged [them] to think about the directionality of my hypothesized relationships and for categorical variables to think

about the effect of each individual category.”

Theme: rTisane provided structure to express and inspect conceptual models.

Participants appreciated how rTisane structured their conceptual modeling process [P2, P4, P9, P10, P11, P12, P13]. Participants found the rTisane DSL particularly helpful. P9 explained how rTisane “led [him] to think about the relationships first, and then whether they were what [he] was hypothesizing” and how this process was the “reverse of the way [he] would think about it normally.” Similarly, P4 explained how using the rTisane language constructs required them to think through how different values of a variable (e.g., different categories) could change income. They observed that their conceptual model with rTisane was “more specific” than without tool support. P4 further explained how rTisane’s DSL “encouraged [them] to think about interactions, which [they] hadn’t thought about before using rTisane.” Four participants said that rTisane generally made it easier for them to specify their conceptual models [P4, P8, P10, P12]. P4 and P10 even believed that rTisane’s “formal structure made [conceptual modeling] more rigorous” [P4] and “more disciplined” [P10].

Participants relied on the conceptual disambiguation step to verify that what they expressed in code accurately represented their implicit assumptions [P2, P8, P12]. P2, who had drawn a conceptual model as a graph on his own prior to using rTisane, said, “The interactive process was a good way to check that the graph came out the same way I was picturing it. It was helpful because it is easier to look at than code” [P2].

Theme: rTisane is expressive enough to capture analysts’ conceptual models accurately.

Importantly, rTisane scaffolded the conceptual modeling process without compromising expressivity. Five participants reported that rTisane had no perceived influence on their conceptual models [P3, P4, P5, P6, P11]. Indeed, three participants expressed identical conceptual models with and without rTisane [P9, P11, P12]. Interestingly, for six participants, the conceptual models they authored with rTisane were subgraphs of conceptual models authored without rTisane [P2, P3, P4, P5, P7, P8]. For P2, P3, P4, and P8, who used rTisane second, rTisane appeared to help focus them on a set of variables and relationships to analyze. P3 explained, “As I started working with rTisane, my first instinct was still to go back to the canvas and do a brainstorming. The process of listing down the categories and the generic relationship between the variables (which was biased to my personal opinion) was still the same (with or without rTisane).” For P5 and P7, who used rTisane first, rTisane provided a starting conceptual model expand upon on their own. For example, P7 authored a statistical model involving an interaction between variables in their rTisane conceptual model when asked to specify a conceptual model on their own. It seems that just conceptual

modeling with rTisane helped P7 translate a conceptual model to a statistical model on his own. Taking these observations together, it seems that rTisane’s DSL can support both convergent and divergent creative thinking about analysts’ domain knowledge.

RQ2: rTisane’s Impact on Statistical Models

Key takeaway: With the exception of picking family and link functions, rTisane focused participants on their analysis goals over low-level details that bogged them down without tool support. As a result, rTisane improved the statistical model authoring process, output statistical models, and communication about statistical analyses.

On their own, three participants were not able to author a statistical model due to unfamiliarity with statistical methods [P3], lack of time [P5], and reliance on visual analyses (ie.g., heatmaps, scatterplots) [P12]. Of the remaining nine participants who completed the study, six participants successfully authored linear regression models [P2, P4, P7, P8, P9, P10]. A seventh participant, P6, started to author a logistic regression model with Race and Income but stopped before binarizing either variable. Two participants, both of whom had just finished authoring statistical models with rTisane, implemented GLMs [P11, P13]. Despite task instructions, P11 started from the rTisane output model script to author their own. After observing the model’s “AIC is large, the residual is large” P11 determined “I don’t think this [rTisane output model] is the right fit.” So, they log-transformed the income variable and fit a new statistical model. P11’s experience mirrors how we anticipate analysts to build upon rTisane output statistical models in the future.

Theme: Statistical models authored with rTisane fit the data just as well or better than statistical models without rTisane.

Of the eight participants who successfully authored linear regression or generalized linear models on their own, three implemented identical models with or without rTisane [P7, P9, P13]. Notably, all three had authored the statistical model with rTisane first, suggesting that rTisane may have biased their own modeling process. For another three participants [P4, P8, P10], their statistical models with rTisane had lower AIC and BIC scores than the statistical models without rTisane. In other words, rTisane models fit the data better or equally well for six out of eight participants. For P11, the statistical model they authored without rTisane dropped some observations, so the models are not directly comparable. For P2, the rTisane statistical model fit worse than his own statistical model in part due to an observed change in his motivation for analysis, discussed below.

Theme: Without rTisane, analysts change their analysis intent during statistical modeling.

Without rTisane, participants [P2, P5, P6, P8, P10], adopted a more exploratory or “data-driven” approach, changing their analysis goals while authoring statistical models. This theme is best illustrated by P2, who started with a hypothesis that Occupation, or Employment, influenced Income. His conceptual model in rTisane had the variables Education, Age, Race, and Sex causing Occupation, which in turn, causes Income (Figure 5.5).

He started authoring statistical models with the intent to assess this hypothesis. On his own, he first authored an ANOVA with Occupation as the IV and Income as the DV. Once he saw that Occupation had a statistically significant influence on Income, he changed his analysis goal to assessing if the variables causing Occupation would “be able to predict which occupation...And then...the income from the occupation just because that’s how I like structured it [in the conceptual model] initially.” However, P2 got stuck on how to author a model with Occupation as the outcome variable because it was categorical, saying, “But the way I structured it in like the diagram. I’m not sure exactly how to do that, because Occupation’s like categorical. Um, so I’m not sure like exactly...how to model that.” This roadblock led P2 to consider an alternative “regression model with Income as like the output and then...all [the IVs] as terms and then just include the interactions between Occupation and the terms that were pointing into it, and that would just be one model.” In other words, P2 tried to author a single statistical model to assess if there was evidence for his conceptual model. However, he was unaware of three key things. First, given his conceptual model, he did not need to account for the other variables to estimate the influence of Occupation on Income and assess his hypothesis. Second, adding interaction terms would not capture the dependencies in the conceptual model. Third, what P2 likely needs to assess all the relationships in his conceptual model is a structural equation model.

While it is well documented that statistical analysis is an iterative process [GW14; JBDM⁺22a] and we saw evidence of this among participants [P5, P6, P10, P11, P12], what P2’s experience exemplifies is how creative participants can be in convincing themselves that the statistical model they authored not only assessed a particular hypothesis but could also arbitrate if their conceptual models were supported by data. Furthermore, this suggests an opportunity for rTisane to support a more iterative analysis process and help analysts author multiple models to assess an entire conceptual model, not just the influence of a single independent variable on a dependent variable.

Theme: Without rTisane, analysts find statistical model formulation challenging.

Participants reported formulating and evaluating statistical models [P2, P3, P5, P8, P12], programming [P6, P13], and preparing data [P7] as the major challenges to authoring statistical models

without rTisane. For example, P3 explained how “There are a number of statistical tests and it gets confusing if I don’t practice it frequently. This is what happened today, I haven’t worked on a hypothesis testing problem recently and while I knew what libraries to go to, I was not sure which test to implement.” Similarly, discussing the details of which covariates to include in a statistical model given a conceptual model, P8 explained how he was uncertain about which “upstream relationships,” or indirect causes, to include in a statistical model. Without rTisane, he described statistical model authoring as “It immediately feels harder doing it directly [without rTisane] like this” [P8].

Theme: rTisane focused analysts on their motivation for analysis.

In contrast, participants reported that rTisane guided them to think about their domains more [P2, P12], lightened their burden in authoring statistical models [P10], and even promoted research transparency [P5] and reproducibility [P4]. Furthermore, rTisane reinforced prior knowledge about statistical methods [P6, P11] and helped participants learn more about GLMs [P4, P6, P7, P13]. P6, who had tried to author a logistic regression model on her own, explained how she could apply what she learned from using rTisane to future analyses: “I like that a multivariate linear regression was used because this will inform any future data analysis...”

Theme: rTisane needs to provide more support for selecting family and link functions.

Despite benefiting from rTisane, many participants had difficulty picking family and link functions in rTisane [P2, P4, P5, P9, P10, P11]. P4 explained, “I didn’t understand the benefit or tradeoffs between different specifications. It wasn’t obvious to me how to create a linear OLS regression, or why I would want to use a specification besides linear OLS.” Given how frequently participants described rTisane as facilitating higher-level of thinking, we attribute the difficulty of selecting family and link functions to the stark contrast between rTisane’s relatively high-level conceptual modeling abstractions and the low-level nature of selecting family and link functions. In the future, it will likely be more usable for rTisane to suggest a specific pair and explain its suggestion rather than require the analyst to pick.

Theme: Analysts want to use rTisane for scientific communication, not just statistical authoring.

When asked how they might imagine using rTisane, participants identified two groups who would benefit: analysts regardless of experience and less technical team members. First, participants described how experienced and novice analysts alike would benefit from using rTisane [P2, P4, P9, P10, P12]. Second, participants mentioned how conceptual models written using rTisane could

be used as boundary objects [SG89] in collaboration with less technical stakeholders [P8, P9]. P8 detailed how a conceptual model written using rTisane could be a communication tool, saying how the “visual representation would play a role in a dialogue with the PI.” P8 went on to imagine how he would like to use rTisane’s conceptual model to generate process diagrams in scientific papers. In other words, how rTisane’s conceptual model could serve as an intermediate representation for multiple kinds of outputs, not just statistical models.

5.8.4 Discussion

rTisane benefits analysts’ conceptual models and statistical models. rTisane’s DSL is expressive to capture analysts’ diverse, nuanced conceptual models. In addition, Importantly, the DSL’s language constructs served as a starting point for statistical analysts to reflect on their domain knowledge. A consequence of rTisane’s DSL and interactive compilation process is that some participants were able to author statistical analyses that they were not able to author on their own. Others could author statistical models that fit the data better than their own statistical models. These results highlight three key insights in rTisane: the benefits of a formalism, balancing usability and rigor, and the potential for re-purposing the intermediate representation.

While unbounded expression in natural language, especially in the era of ChatGPT [BMR⁺20], is enticing, we found that participants found the prospect of expressing their conceptual models using any means daunting. A key benefit of rTisane is that it introduces a formalism that productively reduces the potentially infinite space of how to express conceptual relationships into a finite set that is expressible in the API. Furthermore, based on feedback from participants in the summative evaluation, it seems that the DSL is effective because it is not only expressive but also usable, which we attribute to our iterative language design process involving end-users. Moreover, learning to use rTisane’s formalism required participants to reflect on their domain knowledge. This highlights how a DSL structures the specification and can turn the process of specification into a reflective activity. In this regard, the conceptual disambiguation step was critical. The graph visualization in the GUI helped analysts reflect on what they expressed and how to resolve any ambiguities present.

A key challenge in designing rTisane was balancing usability and rigor. On one hand, we wanted to make it easy for analysts to express their conceptual models (usability), but we also wanted to ensure that the conceptual models they expressed were amenable to formal causal reasoning to derive statistical models (rigor). We were able to achieve both in rTisane by designing usable language constructs in the DSL and increasing precision for rigor during disambiguation.

Finally, participants discussed the potential for using conceptual models to communicate with less technical collaborators. Implicit in this recommendation is an acknowledgement of the usefulness of a conceptual model as an intermediate representation. While rTisane is focused on using

conceptual models to derive statistical models, there may be additional “backends” to target, such as scientific model diagramming or planning study procedures.

5.8.5 Limitations and Future Work on rTisane

While participants found rTisane helpful, they suggested three major areas of improvement: (i) statistical model interpretation, (ii) iterative model revision, (iii) system usability.

Once analysts execute the output statistical model from rTisane, they find the output results too low-level. Because rTisane uses lme4 under the hood, the outputs are the default model outputs from lme4. However, given that rTisane’s input language is at the conceptual level, analysts expected the outputs to at least relate back to the conceptual model they input. In other words, the input and output levels of abstraction should be commensurate. This support would facilitate what analysts already try to do with statistical analyses they author on their own without rTisane. P8 found the output from lme4 overwhelming, saying, “Looking at the summary() in R was too much to look at.” He suggested a simple way to tie the results back to his input conceptual model: “Would be nice if you could have the same visual representation with p-values/coefficients!”

Furthermore, while participants could iterate on their conceptual models by adding or removing variables and relationships, they could not engage in a larger iteration loop with their output statistical model from rTisane. Improving statistical result interpretation would help with model iteration. In addition, participants also sought more direct support. For instance, P11 described the rTisane output statistical model as “an initial or baseline model but follow-up evaluation of the model is needed.” They wanted to “go back and tweak things a bit” about their statistical model. This kind of model iteration is not only typical of the participants’ workflows but also even a best practice recommendation from the statistics community [GVS⁺20b]. Supporting novice and more experienced analysts revise models will likely require different levels of abstraction and automation.

Finally, participants found going back and forth between code and an interface outside their IDE complicated and “clunky.” While part of this may have been in part due to the fact that participants were using rTisane on a remote desktop, embedding rTisane in a notebook seems likely to reduce major usability issues. Additionally, participants gave suggestions for syntactic sugar for specifying conceptual models. For example, instead of specifying multiple causes and relates statements, they wished they could batch specify and add them to the conceptual model. Ways to reduce the specification burden for analysts by providing syntactic sugar or even removing the need to program at all are interesting avenues to explore.

5.9 Discussion

Iteratively designing, developing, and evaluating Tisane’s key insight—to compile statistical models from conceptual models—led to rTisane. A key step in designing rTisane was the exploratory lab study that used Tisane as a probe. The exploratory lab study suggested the need to allow analysts to express their conceptual models using more granular, low-level functions. Although obvious in hindsight, this finding was *counterintuitive* at the time.

A widely held belief, especially within the HCI community, is that the higher the level of abstraction for a task, the better for end-users. However, we saw the opposite. Statistical non-experts engaged deeply with conceptual models about their domain and wanted to be more detailed and specific when describing their conceptual models. In other words, while the focus on the abstraction should be at the conceptual level, within that, analysts want to move up and down the ladder of abstraction. Indeed, in the summative evaluation of rTisane, we saw that analysts made use of all language constructs and reported finding them instrumental to deepening their consideration of domain knowledge. Arguably, the conceptual modeling language constructs also benefited the quality of statistical models output from rTisane. Using rTisane, analysts authored statistical models that fit the data better than their own or authored identical statistical models after using rTisane.

Based on my experience with Tisane and rTisane, I speculate that abstractions can achieve usability and rigor by matching the content-focus of end-users while giving them opportunities to get into the low-level details of specifications. This gives end-users the agency to express themselves more fully, transforming the programming task from strictly a means to an end to a meaningful, reflective activity in itself.

5.10 Summary of Contributions

Tisane embodies the hypothesis central to this dissertation: A DSL for expressing implicit conceptual knowledge and automated reasoning enable statistical non-experts to author valid statistical models. Through an iterative design and evaluation process, we refined core language constructs for expressing conceptual models (**Thesis Challenge 1: Explicating domain knowledge**) and introduced a two-step interactive disambiguation process for compiling conceptual knowledge into statistical analysis code (**Thesis Challenge 2: Representation and reasoning**). rTisane implements these changes. In case studies of Tisane and a controlled lab study of rTisane, we found that the DSL is expressive enough to capture analysts’ conceptual models accurately, eases the burden of making their implicit assumptions explicit, and pushes analysts to think about their domains more deeply. rTisane’s automated reasoning led to statistical models that fit the data just as well as or better than statistical models authored without. rTisane even helped analysts who were not

able to author statistical models on their own get to an output statistical model. Analysts also reported that through the process, they learned about GLMs (**Thesis Challenge 3: Statistical understanding**). Together, these results show how DSLs and automated reasoning together in fact do help statistical non-experts author valid statistical analyses that, arguably, they would not be able to author without such support.

The first release of Tisane was a collaboration with Audrey Seo, Jeffrey Heer, and René Just. The corresponding paper was originally published and presented at ACM CHI 2022 [JSHJ22], where it received a Best Paper Honorable Mention Award. The exploratory design study, second system iteration, and the summative evaluation are in collaboration with Edward Misback, Jeffrey Heer, and René Just. The corresponding paper is under submission and has not yet been published.

Chapter 6

Conclusion

While statistical analysis has become more pervasive among end-users who are not statistical experts, the tools for conducting analyses have continued to require high statistical expertise. This dissertation examines how to design and develop tools that not only lower the barriers for statistical non-experts but also provide guarantees about the validity of authored analyses. We introduce two new systems, Tea and Tisane. Both provide DSLs for expressing implicit conceptual knowledge and then compile these high-level specifications into statistical analyses, Null Hypothesis Significance Tests in Tea (Chapter 3) and generalized linear models with or without mixed effects in Tisane (Chapter 5). Additionally, we develop a theory of hypothesis formalization that describes the cognitive and operational steps involved in translating a conceptual research question into a statistical analysis implementation in code. Our theory of hypothesis formalization retrospectively validated our design in Tea and directly inspired the design of Tisane.

Relate this work to chasm/bridge that Amelia identified in her work Existing statistical analysis tools are either designed for students learning basic statistics or statistical experts [McN15]. Tools do not support statistical non-experts, such as researchers, through the process of authoring accurate statistical models.

6.1 Discussion

This dissertation addresses three challenges central to the thesis that (i) programming abstractions focused on capturing analysts' implicit conceptual knowledge and (ii) formal representations and reasoning to determine statistical analyses benefit statistical non-experts. We discuss each challenge and how the projects in this dissertation address them below.

6.1.1 Challenge 1: Designing the *right* level of abstraction

With any programming language, end-users must learn and use a formalism. Tea and Tisane provide high-level abstractions but the key to their design is that they abstract the appropriate *conceptual concerns* implicitly involved in statistical analyses. In fact, the fact that an abstraction is high or low is less relevant. Indeed, a key insight that guided our design of rTisane (Section 5.7) was that analysts wanted low- and high-level conceptual abstractions to express their domain knowledge with varying degrees of detail that felt helpful and accurate to them (see Section 5.6).

When comparing the abstractions Tea and Tisane provide, it is easy to see that the conceptual relationships between variables were still largely implicit in Tea. An important takeaway from the theory of *hypothesis formalization* was the importance of conceptual models, which are present for statistical testing and modeling alike. Therefore, conceptual models should be a central concern in designing programming abstractions for data analysis.

6.1.2 Challenge 2: Representing and reasoning about analysis decisions

The abstractions that may be usable to statistical non-experts may not be precise enough for formal reasoning (Section 5.6). Therefore, a key challenge in designing representations amenable to reasoning is in finding a “shared representation” between analysts and computational techniques. Based on Tea’s key insight that statistical test selection can be reformulated as a constraint satisfaction problem, we represented statistical tests using logical constraints in a knowledge base. Using Tea’s DSL, analysts specify additional constraints about their hypothesis and data, which helps Tea’s runtime system solve a system of constraints to identify valid statistical tests. In Tisane, the shared representation is the conceptual model, which Tisane represents as a graph. This representation made reasoning about linear model formulations straightforward by applying causal reasoning techniques on a part of the graph.

In designing these shared representations, a temptation was to fit the DSL on top of a reasoning approach that was straightforward. In this view, the DSL would be a thin wrapper around the automated reasoning engine. For example, a very early prototype of Tisane used Answer Set Programming (ASP) to define when specific confounders should appear in a generalized linear model. In addition to being a clunky way to represent linear model formulation rules when the statistics community has converged on using graphs, this prototype required analysts to incrementally refine their statistical models by interacting with the UNSAT core. This interaction model, though interesting, did not allow us to discovery and fully realize the real benefit of expressing conceptual models: giving analysts an opportunity to reflect on their assumptions in an open-ended way.

6.1.3 Challenge 3: Interaction as reflection

As we saw in the case studies with Tisane, providing abstractions and interactions with shared representations for formal reasoning increases analysts' awareness of their implicit assumptions, data, and analysis practices. By providing the appropriate abstractions, DSLs can make the specification process a useful form of documentation. This may later be useful for sharing and inspection. For instance, by stating their implicit conceptual and data assumptions in Tea and Tisane, researchers can help improve scientific replicability and reproducibility.

6.2 Recent developments

Mention: in LLMs impact the contributions of this dissertation, exciting opportunities to leverage them to realize the goals of this work

6.2.1 Construct validity: Within reach with the usage of LLMs

This thesis focused on internal, external, and statistical conclusion validity. However, could reason about construct validity with LLMs.

6.2.2 What about in the face of LLMs?

But how do people express their domain knowledge, make the process meaningful

Mention LLMs as a technology to use here?

6.3 Limitations and Future work

To do (??)

This dissertation scrutinizes how statistical non-experts author statistical analyses. We elaborate one some of the limitations of this work and opportunities for future research.

6.3.1 Support interpretation of statistical results

While Tisane effectively addresses the gulf of execution by compiling conceptual models into statistical models, it falls short of bridging the gulf of evaluation. Tisane does not yet provide support for analysts to interpret the results of their statistical models. Future research should focus on two related challenges: (i) improving statistical reporting to enhance the understanding of results and (ii) providing support for navigating the consequences of the results, such as updates to analysts' conceptual models or the need to revise statistical models.

6.3.2 Connect statistical modeling and testing

A natural question arising from this dissertation is the choice between Tea and Tisane for analysts. Tea focuses on statistical testing, determining evidence for or against a specific claim, while Tisane emphasizes statistical modeling, estimating variable influences in the presence of other variables (e.g., confounders, mediators, etc.) However, statistical testing and modeling are not mutually exclusive. Analysts often want to conduct tests after building models to answer substantive questions as well as assess model fit. A compelling future direction is to enable analysts to ask follow-up questions about specific estimates and effects from a statistical model.

6.3.3 Develop a grammar of study design

There are relationships between Tisane’s language constructs for specifying data collection details, data schema specification, and experimental design. How could we draw these connections out and formalize them?

Tisane’s graph IR is an entity-relationship (ER) model [Che76]. ER models, or diagrams, are used to describe data schema. ER models describe how entities relate to other entities and attributes. In Tisane, a variable of Unit type can be viewed as an entity. The `nests_within` edge describes how two units, or entities, relate to one another. Tisane’s graph IR also relates units (entities) to measures (attributes).

6.3.4 Support more phases of the data lifecycle

This dissertation emphasizes the need for abstractions that capture analysts’ implicit domain knowledge. These abstractions enable valid analysis formulation and promote reflective thinking. Building upon this, we can begin to ask how the same ideas—abstractions and automated reasoning for conceptual knowledge often implicit in statistical analyses—could apply to other phases of the data lifecycle. Future work should explore how to elicit and track the evolution of conceptual knowledge even before statistical analysis by developing new elicitation techniques and representations of domain knowledge and ecosystem of inter-operating tools to track and ensure validity throughout the data lifecycle.

6.3.5 Promote analytical best practices in science

Tea and Tisane primarily follow a top-down authoring approach, where analysts start with a research question and hypothesis. However, as observed in our lab study to develop hypothesis formalization (Section 4.3), analysts often develop and refine hypotheses based on statistical results. Therefore, a future research direction would be to develop ways to incorporate both data-driven and research

question-driven approaches to model authoring and refinement that do not promote cherry-picking. One possibility is to leverage Tea’s and Tisane’s reasoning capabilities to reason in multiple other directions, from statistical models and data to all possible conceptual models or statistical models to data invariants that could inform study designs.

Moreover, one of the precautions integrated into the design of Tisane was to prevent cherry-picking and p-hacking by using analysts’ conceptual models to drive statistical model formulation. While Tisane supports mapping one conceptual model to a statistical model, an under-explored direction is to assess the robustness of effects across multiple possible conceptual models, especially in cases of ambiguity or debate in a discipline. On the other hand, multiverse analyses [] embrace conceptual model uncertainty by considering all possible statistical model formulations. Future research should look into how to both be consistent with aspects of conceptual models researchers know and assess evidence for other competing aspects of conceptual models.

6.3.6 Improving data science education

1 sentence: Motivation: From conversations with students and my own experience taking statistics courses throughout undergraduate and graduate education is that the connection between statistical methods and the kinds of questions I want to ask is often unclear. Furthermore, students feel they need to memorize a bunch of different methods, at the cost of thinking through what their substantive questions are and what they should care about in a statistical approach.

2: Goal, Impact: Promote statistical thinking, an essential component of practicing data science.

Especially, greatest potential: Especially to teach students to separate multiple sets of concerns, specifically the conceptual from the statsitical

3. Role of systems: Use Tea and Tisane to focus students to focus first on identifying and articulating their motivation and intents for analysis

Before introducing and teaching statistical details

4. Future work: Deploy Tea and Tisane- questions to answer: impact on statistical thinking, computational thinking - what additional tools are necessary to develop?

A ramp from novice to expert tools is missing in the current ecosystem [McN18]. Tea and Tisane lay the foundation for a bridge between novice and expert tools by providing abstractions that match those of statistical non-experts while also giving experts control, flexibility, and compatibility with other expert tools in Python and R. While pursuing the above research agenda, I look forward to directly improving data science courses I teach by deploying my systems in the classroom, discovering students’ needs, and iterating on tools and curriculum. In pursuing this goal, I want to ease transitions between analysis paradigms (e.g., NHST and Bayesian inference). One promising direction is to separate concerns about model specification from interpretation and assessment.

Get rid of?: In a small way, this separation also teaches students how to separate specification from implementation, an essential perspective in computer science. So that if students need to implement more complex analyses, they have some awareness of how to organize their computational approach.

6.4 Impact

Include details of where this code is available, open source, etc. The most rewarding part of conducting the work in this dissertation has been to see real-world use cases and adoption of tools. As of May 2023, Tea has been downloaded 15K times, and the first release of Tisane has been downloaded 12K times. Over the last few years, I have also enjoyed reading and answering a flurry of emails where analysts, including scientists and social scientists, share anecdotes of how they have used (and sometimes failed to use) Tea and Tisane.

Add IHME anecdote/story/example?

To more systematically capture and act on these kinds of anecdotes in the future, I am excited to develop a web platform for Tea and Tisane, where users can share their programs, data, and insights. Over time, I hope to collect a gallery of examples to answer questions about challenges using the DSLs, the learn which use cases are under-supported, and assess the practical impact of using conceptual abstractions and automated reasoning. I would also like to see if end-users repeatedly use Tea and Tisane or if these are one-off engagements. I look forward to not only pursuing ideas and systems developed in this dissertation but also fostering a community of users.

6.5 Closing Remarks

designed and implemented two systems, Tea [JDR⁺19] and Tisane [JSHJ22], that leverage **domain-specific languages** (DSLs) to capture analysts' implicit assumptions and conceptual knowledge. Users **interactively compile** these high-level specifications into low-level code. To infer valid statistical analyses, the systems **programmatically represent and reason about core statistical authoring challenges** as constraints and graphs. As a result, my systems prevent common analysis mistakes [JDR⁺19; JSHJ22].

This thesis furthers our understanding of what makes statistical analyses difficult to author and then designs and implements two domain-specific languages (DSLs) to addressing these issues. The DSLs illustrate that automating aspects of statistical analysis benefits statistical non-experts. However, more importantly, this dissertation illustrates how we can - designing abstractions that capture intent - reifying the connection between domain knowledge and statistical analysis - design-

ing interfaces and interactions that increase analyst awareness of the impact of analysis, rationale.
← change this based on rTisane eval?

Bibliography

- [A⁺83] American Psychological Association et al. *Publication manual*. American Psychological Association Washington, DC, 1983.
- [ACM⁺19] Ricard Argelaguet, Stephen J Clark, Hisham Mohammed, L Carine Stapel, Christel Krueger, Chantriolnt-Andreas Kapourani, Ivan Imaz-Rosshandler, Tim Lohoff, Yunlong Xiang, Courtney W Hanna, et al. Multi-omics profiling of mouse gastrulation at single-cell resolution. *Nature*, 576(7787):487–491, 2019.
- [ADP19] Reena Aggarwal, Sandeep Dahiya, and Nagpurnanand R Prabhala. The power of shareholder votes: Evidence from uncontested director elections. *Journal of Financial Economics*, 133(1):134–153, 2019.
- [Ass96] American Psychological Association. Task force on statistical inference, 1996.
- [AWT20] Ruben C Arslan, Matthias P Walther, and Cyril S Tata. formr: A study framework allowing for automated feedback generation and complex longitudinal experience-sampling studies using r. *Behavior Research Methods*, 52(1):376–387, 2020.
- [AZL⁺18] Sara Alspaugh, Nava Zokaei, Andrea Liu, Cindy Jin, and Marti A Hearst. Futzng and moseying: Interviews with professional data analysts on exploration practices. *IEEE transactions on visualization and computer graphics*, 25(1):22–31, 2018.
- [B⁺17] Paul-Christian Bürkner et al. brms: An r package for bayesian multilevel models using stan. *Journal of statistical software*, 80(1):1–28, 2017.
- [Bak16] Monya Baker. 1,500 scientists lift the lid on reproducibility. *Nature*, 533(7604), 2016.
- [Bar13] Dale J Barr. Random effects structure for testing interactions in linear mixed-effects models. *Frontiers in psychology*, 4:328, 2013.
- [BB16] Paul-Christian Bürkner and Maintainer Paul-Christian Buerkner. Package ‘brms’. 2016.
- [BBC⁺09] Benjamin M Bolker, Mollie E Brooks, Connie J Clark, Shane W Geange, John R Poulsen, M Henry H Stevens, and Jada-Simone S White. Generalized linear mixed models: a practical guide for ecology and evolution. *Trends in ecology & evolution*, 24(3):127–135, 2009.

- [BCCH19a] Graeme Blair, Jasper Cooper, Alexander Coppock, and Macartan Humphreys. Declaring and diagnosing research designs. *American Political Science Review*, 113(3):838–859, 2019.
- [BCCH19b] Graeme Blair, Jasper Cooper, Alexander Coppock, and Macartan Humphreys. Declaring and diagnosing research designs. *American Political Science Review*, 113(3):838–859, 2019.
- [BCJ⁺19] Eli Bingham, Jonathan P Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D Goodman. Pyro: Deep universal probabilistic programming. *The Journal of Machine Learning Research*, 20(1):973–978, 2019.
- [BCK⁺19] Ambre M Bertholet, Edward T Chouchani, Lawrence Kazak, Alessia Angelin, Andriy Fedorenko, Jonathan Z Long, Sara Vidoni, Ryan Garrity, Joonseok Cho, Naohiro Terada, et al. H⁺ transport is an integral function of the mitochondrial adp/atp carrier. *Nature*, 571(7766):515–520, 2019.
- [BCLW18] Leo Breiman, Adele Cutler, Andy Liaw, and Matthew Wiener. Package ‘randomforest’. 2018.
- [BDSDL⁺19] Alexandre P Bénéchet, Giorgia De Simone, Pietro Di Lucia, Francesco Cilenti, Giulia Barbiera, Nina Le Bert, Valeria Fumagalli, Eleonora Lusito, Federica Moalli, Valentina Bianchessi, et al. Dynamics and genomic landscape of cd8+ t cells undergoing hepatic priming. *Nature*, 574(7777):200–205, 2019.
- [BEB14] Eytan Bakshy, Dean Eckles, and Michael S Bernstein. Designing and deploying online field experiments. In *Proceedings of the 23rd international conference on World wide web*, pages 283–292. ACM, 2014.
- [Ber11] Jacques Bertin. *Graphics and graphic information processing*. Walter de Gruyter, 2011.
- [Bet20] Michael Betancourt. Towards a principled bayesian workflow, 2020.
- [BH19] Leilani Battle and Jeffrey Heer. Characterizing exploratory visual analysis: A literature review and evaluation of analytic provenance in tableau. *Computer Graphics Forum (Proc. EuroVis)*, 2019.
- [BHM19] Farah Bughio, Gary R. Huckell, and Keith A. Maggert. Monitoring of switches in heterochromatin-induced silencing shows incomplete establishment and developmental instabilities. 116(40):20043–20053, 2019.
- [BKvB⁺17] Mollie E Brooks, Kasper Kristensen, Koen J van Benthem, Arni Magnusson, Casper W Berg, Anders Nielsen, Hans J Skaug, Martin Machler, and Benjamin M Bolker. glmmTMB balances speed and flexibility among packages for zero-inflated generalized linear mixed modeling. *The R journal*, 9(2):378–400, 2017.
- [BLB⁺13] Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort,

- Jaques Grobler, et al. Api design for machine learning software: experiences from the scikit-learn project. *arXiv preprint arXiv:1309.0238*, 2013.
- [BLST13] Dale J Barr, Roger Levy, Christoph Scheepers, and Harry J Tily. Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of memory and language*, 68(3):255–278, 2013.
- [BLT⁺19] Natalie Biderman, Roy Luria, Andrei R. Teodorescu, Ron Hajaj, and Yonatan Goshen-Gottstein. Working memory has better fidelity than long-term memory: The fidelity constraint is not a general property of memory after all. *Psychological Science*, 30(2):223–237, 2019.
- [BMB⁺19] Douglas Bates, Martin Mächler, Ben Bolker, Steve Walker, Rune H.B. Christensen, Henrik Singmann, Bin Dai, Fabian Scheipl, Gabor Grothendieck, Peter Green, and John Fox. Package ‘lme4’. *CRAN*, 2019.
- [BMBW14] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. Fitting linear mixed-effects models using lme4. *arXiv preprint arXiv:1406.5823*, 2014.
- [BMR⁺20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [Bob17] Igor Bobriakov. Top 15 python libraries for data science in 2017. *ActiveWizards in Medium*, 2017.
- [Bob18] Igor Bobriakov. Top 20 python libraries for data science in 2018. *ActiveWizards in Medium*, 2018.
- [Bru19] J. Bruin. Choosing the correct statistical test in sas, stata, spss and r, 2019.
- [BSCN19] Timothy Ballard, David K. Sewell, Daniel Cosgrove, and Andrew Neal. Information processing under reward versus under punishment. *Psychological Science*, 30(5):757–764, 2019.
- [BSG12] Mauro Borgo, Alessandro Soranzo, and Massimo Grassi. Psychtoolbox: sound, keyboard and mouse. In *MATLAB for Psychologists*, pages 249–273. Springer, 2012.
- [C⁺15] François Chollet et al. Keras. <https://keras.io>, 2015.
- [Cai07] Paul Cairns. Hci... not as it should be: inferential statistics in hci research. In *Proceedings of the 21st British HCI Group Annual Conference on People and Computers: HCI... but not as we know it- Volume 1*, pages 195–201. British Computer Society, 2007.
- [CCM19a] Andrea Caggese, Vicente Cuñat, and Daniel Metzger. Firing the wrong workers: Financing constraints and labor misallocation. *Journal of Financial Economics*, 133(3):589–607, 2019.

- [CCM19b] Andrea Caggese, Vicente Cuñat, and Daniel Metzger. Firing the wrong workers: Financing constraints and labor misallocation. *Journal of Financial Economics*, 133(3):589–607, 2019.
- [CCWA13] Jacob Cohen, Patricia Cohen, Stephen G West, and Leona S Aiken. *Applied multiple regression/correlation analysis for the behavioral sciences*. Routledge, 2013.
- [CDdME19] Samuel H Church, Seth Donoughe, Bruno AS de Medeiros, and Cassandra G Extavour. Insect egg size and shape evolve with ecology but not developmental rate. *Nature*, 571(7763):58–62, 2019.
- [CEG⁺16] Robert Carver, Michelle Everson, John Gabrosek, Nicholas Horton, Robin Lock, Megan Mocko, Allan Rossman, Ginger Holmes Roswell, Paul Velleman, Jeffrey Witmer, et al. Guidelines for assessment and instruction in statistics education (gaise) college report 2016. 2016.
- [CFP20] Carlos Cinelli, Andrew Forney, and Judea Pearl. A crash course in good and bad controls. *Sociological Methods & Research*, page 00491241221099552, 2020.
- [CGDH⁺17] Bob Carpenter, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan : A probabilistic programming language. *Journal of Statistical Software*, 76, 01 2017.
- [Che76] Peter Pin-Shan Chen. The entity-relationship model—toward a unified view of data. *ACM transactions on database systems (TODS)*, 1(1):9–36, 1976.
- [CKM19] Alan D Crane, Andrew Koch, and Sébastien Michenaud. Institutional investor cliques and governance. *Journal of Financial Economics*, 133(1):175–197, 2019.
- [Cla73] Herbert H Clark. The language-as-fixed-effect fallacy: A critique of language statistics in psychological research. *Journal of verbal learning and verbal behavior*, 12(4):335–359, 1973.
- [CLD⁺19] Nadia Chernyak, Kristin L. Leimgruber, Yarrow C. Dunham, Jingshi Hu, and Peter R. Blake. Paying back people who harmed us but not people who helped us: Direct negative reciprocity precedes direct positive reciprocity in early development. *Psychological Science*, 30(9):1273–1286, 2019.
- [CLM⁺20] Yunshun Chen, Aaron TL Lun, Davis J McCarthy, Matthew E Ritchie, Belinda Phipson, Yifang Hu, Xiaobei Zhou, Mark D Robinson, and Gordon K Smyth. Empirical analysis of digital gene expression data in r (v3.30.3). 2020.
- [CMB18] Sarah E Chasins, Maria Mueller, and Rastislav Bodik. Rousillon: Scraping distributed hierarchical web data. In *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*, pages 963–975, 2018.
- [CMR⁺19] Meng S. Choy, Thomas M. Moon, Rini Ravindran, Johnny A. Bray, Lucy C. Robinson, Tara L. Archuleta, Wuxian Shi, Wolfgang Peti, Kelly Tatchell, and Rebecca Page. Sds22 selectively recognizes and traps metal-deficient inactive pp1. 116(41):20472–20481, 2019.

- [CMR⁺20] Yunshun Chen, David McCarthy, Matthew Ritchie, Mark Robinson, and Gordon Smyth. edger: differential analysis of sequence read count data. 2020.
- [Coh88] Jacob Cohen. Statistical power analysis for the social sciences. 1988.
- [Com23] Plotly Dash Community. Plotly dash, 2023.
- [Cru19] Matthew J.C. Crump. Jspysch, 2019.
- [CWM⁺19] Ethan C Campbell, Earle A Wilson, GW Kent Moore, Stephen C Riser, Casey E Brayton, Matthew R Mazloff, and Lynne D Talley. Antarctic offshore polynyas linked to southern hemisphere climate anomalies. *Nature*, 570(7761):319–325, 2019.
- [DL15a] Joshua R De Leeuw. jspsych: A javascript library for creating behavioral experiments in a web browser. *Behavior research methods*, 47(1):1–12, 2015.
- [DL15b] Joshua R De Leeuw. jspsych: A javascript library for creating behavioral experiments in a web browser. *Behavior research methods*, 47(1):1–12, 2015.
- [DMB08] Leonardo De Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340. Springer, 2008.
- [DMC⁺19] Florent Détröit, Armand Salvador Mijares, Julien Corny, Guillaume Daver, Clément Zanolli, Eusebio Dizon, Emil Robles, Rainer Grün, and Philip J Piper. A new species of homo from the late pleistocene of the philippines. *Nature*, 568(7751):181–186, 2019.
- [DMFKS19] Marco Di Maggio, Francesco Franzoni, Amir Kermani, and Carlo Sommavilla. The relevance of broker networks for information diffusion in the stock market. *Journal of Financial Economics*, 134(2):419–446, 2019.
- [Efr92] Bradley Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics*, pages 569–593. Springer, 1992.
- [Ehr73] Isaac Ehrlich. Participation in illegitimate activities: A theoretical and empirical investigation. *Journal of political Economy*, 81(3):521–565, 1973.
- [EVL⁺19] Motahhare Eslami, Kristen Vaccaro, Min Kyung Lee, Amit Elazari Bar On, Eric Gilbert, and Karrie Karahalios. User attitudes towards algorithmic opacity and transparency in online reviewing platforms. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2019.
- [EWBLM19] Alexander Eiselmayer, Chatchanan Wacharamanotham, Michel Beaudouin-Lafon, and Wendy Mackay. Touchstone2: An interactive environment for exploring trade-offs in hci experiment design. 2019.
- [FHT⁺20] Jerome Friedman, Trevor Hastie, Rob Tibshirani, Balasubramanian Narasimhan, Kenneth Tay, Noah Simon, and Junyang Qian. Package ‘glmnet’. 2020.
- [Fis37] Ronald Aylmer Fisher. *The design of experiments*. Oliver And Boyd; Edinburgh; London, 1937.

- [FKE⁺15] Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [FKL⁺12] Janet Feigenspan, Christian Kästner, Jörg Liebig, Sven Apel, and Stefan Hanenberg. Measuring programming experience. In *2012 20th IEEE International Conference on Program Comprehension (ICPC)*, pages 73–82. IEEE, 2012.
- [FMF12] Andy Field, Jeremy Miles, and Zoë Field. *Discovering statistics using R*. Sage publications, 2012.
- [FRG⁺19] Janos Fuzik, Sabah Rehman, Fatima Girach, Andras G Miklosi, Solomiia Korchynska, Gloria Arque, Roman A Romanov, János Hanics, Ludwig Wagner, Konstantinos Meletis, et al. Brain-wide genetic mapping identifies the indusium griseum as a prenatal target of pharmacologically unrelated psychostimulants. *Proceedings of the National Academy of Sciences*, 116(51):25958–25967, 2019.
- [FSD⁺19] Dylan Fafard, Ian Stavness, Martin Dechant, Regan Mandryk, Qian Zhou, and Sidney Fels. Ftvr in vr: Evaluation of 3d perception with a simulated volumetric fish-tank virtual reality display. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2019.
- [GAL⁺19] Lukas Gehrke, Sezen Akman, Pedro Lopes, Albert Chen, Avinash Kumar Singh, Hsiang-Ting Chen, Chin-Teng Lin, and Klaus Gramann. Detecting visuo-haptic mismatches in virtual reality using the prediction error negativity of event-related brain potentials. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–11, 2019.
- [GCS⁺13] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC press, 2013.
- [GDH07] François Guimbretière, Morgan Dixon, and Ken Hinckley. Experiscope: an analysis tool for interaction data. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 1333–1342. ACM, 2007.
- [Gel05] Andrew Gelman. Why i don’t use the term “fixed and random effects”, 2005.
- [GH06] Andrew Gelman and Jennifer Hill. *Data analysis using regression and multi-level/hierarchical models*. Cambridge university press, 2006.
- [GKM19] John M Griffin, Samuel Kruger, and Gonzalo Maturana. Do labor markets discipline? evidence from rmbs bankers. *Journal of Financial Economics*, 133(3):726–750, 2019.
- [GL13] Andrew Gelman and Eric Loken. The garden of forking paths: Why multiple comparisons can be a problem, even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited ahead of time. *Department of Statistics, Columbia University*, 2013.

- [GPR99] Sander Greenland, Judea Pearl, and James M Robins. Causal diagrams for epidemiologic research. *Epidemiology*, pages 37–48, 1999.
- [Gro19] Garrett Grolemund. Quick list of useful r packages. *R Studio Support*, 2019.
- [GS20] LLC. GraphPad Software. Graphpad prism 8 user guide. 2020.
- [GSV⁺19] Jonah Gabry, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman. Visualization in bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(2):389–402, 2019.
- [GSY19] Robin Greenwood, Andrei Shleifer, and Yang You. Bubbles for fama. *Journal of Financial Economics*, 131(1):20–43, 2019.
- [GVS⁺20a] Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C Margossian, Bob Carpenter, Yuling Yao, Lauren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. Bayesian workflow. *arXiv preprint arXiv:2011.01808*, 2020.
- [GVS⁺20b] Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C Margossian, Bob Carpenter, Yuling Yao, Lauren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. Bayesian workflow. *arXiv preprint arXiv:2011.01808*, 2020.
- [GW14] Garrett Grolemund and Hadley Wickham. A cognitive interpretation of data analysis. *International Statistical Review*, 82(2):184–204, 2014.
- [H⁺10] Jarrod D Hadfield et al. Mcmc methods for multi-response generalized linear mixed models: the mcmcglmm r package. *Journal of Statistical Software*, 33(2):1–22, 2010.
- [Had20] Jarrod Hadfield. Package ‘mcmcglmm’. 2020.
- [HAPGNN⁺19] Mireia Hernández, María Ángeles Palomar-García, Benito Nohales-Nieto, Gustau Olcina-Sempere, Esteban Villar-Rodríguez, Raúl Pastor, César Ávila, and Maria-Antónia Parcet. Separate contribution of striatum volume and pitch discrimination to individual differences in music reward. *Psychological Science*, 30(9):1352–1361, 2019.
- [Har20] Peter Harrison. psychtestr: An r package for designing and conducting behavioural psychological experiments. *The Journal of Open Source Software*, 5(49), 2020.
- [HBH18] Jane Hoffswell, Alan Borning, and Jeffrey Heer. Setcola: High-level constraints for graph layout. In *Computer Graphics Forum*, volume 37, pages 537–548. Wiley Online Library, 2018.
- [HHNT89] John H Holland, Keith J Holyoak, Richard E Nisbett, and Paul R Thagard. *Induction: Processes of inference, learning, and discovery*. MIT press, 1989.
- [HHZ19] Pinghsun Huang, Hsin-Yi Huang, and Yan Zhang. Do firms hedge with foreign currency derivatives for employees? *Journal of Financial Economics*, 133(2):418–440, 2019.

- [HLC19] Brian Hempel, Justin Lubin, and Ravi Chugh. Sketch-n-sketch: Output-directed programming for svg. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*, pages 281–292, 2019.
- [HM19] Sebok K. Halder and Richard Milner. A critical role for microglia in maintaining vascular integrity in the hypoxic spinal cord. 116(51):26029–26037, 2019.
- [Hol79] Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics*, pages 65–70, 1979.
- [Hoy95] Rick H Hoyle. *Structural equation modeling: Concepts, issues, and applications*. Sage, 1995.
- [HQ14] Trevor Hastie and Junyang Qian. Glmnet vignette. 2014.
- [HWV⁺19] Nur Al-huda Hamdan, Adrian Wagner, Simon Voelker, Jürgen Steimle, and Jan Borchers. Springlets: Expressive, flexible and silent on-skin tactile interfaces. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2019.
- [JBDM⁺22a] Eunice Jun, Melissa Birchfield, Nicole De Moura, Jeffrey Heer, and Rene Just. Hypothesis formalization: Empirical findings, software limitations, and design implications. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 29(1):1–28, 2022.
- [JBDM⁺22b] Eunice Jun, Melissa Birchfield, Nicole De Moura, Jeffrey Heer, and Rene Just. Hypothesis formalization: Empirical findings, software limitations, and design implications. In *ACM Transactions on Computer-Human Interaction (TOCHI)*, volume 29, 2022.
- [JDR⁺19] Eunice Jun, Maureen Daum, Jared Roesch, Sarah E Chasins, Emery D Berger, Rene Just, and Katharina Reinecke. Tea: A high-level language and runtime system for automating statistical analysis. In *Proceedings of the 32nd Annual Symposium on User Interface Software and Technology*. ACM, 2019.
- [JMHJ23] Eunice Jun, Edward Misback, Jeffrey Heer, and René Just. rtisane: Formalizing conceptual models to author statistical models reduces error, increases awareness, and teaches novices. In *Under submission*, 2023.
- [Jol18] Eshin Jolly. Pymer4: connecting r and python for linear mixed modeling. *Journal of Open Source Software*, 3(31):862, 2018.
- [JOP⁺21a] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–2021.
- [JOP⁺21b] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–2021.
- [JOP⁺21c] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–2021.

- [JS11] Bradley Jones and John Sall. Jmp statistical discovery software. *Wiley Interdisciplinary Reviews: Computational Statistics*, 3(3):188–194, 2011.
- [JS23] Eunice Jun and Audrey Seo. Tisane, 2022-2023.
- [JSHJ22] Eunice Jun, Audrey Seo, Jeffrey Heer, and René Just. Tisane: Authoring statistical models via formal reasoning from conceptual and data relationships. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–16, 2022.
- [JSP⁺19] Robert J Johnston, Linhui Julie Su, Jason Pinckney, David Critton, Eric Boyer, Arathi Krishnakumar, Martin Corbett, Andrew L Rankin, Rose Dibella, Lynne Campbell, et al. Vista is an acidic ph-selective ligand for psgl-1. *Nature*, 574(7779):565–570, 2019.
- [Kab11] Robert I Kabacoff. R: In action. 2011.
- [KD88] David Klahr and Kevin Dunbar. Dual space search during scientific reasoning. *Cognitive science*, 12(1):1–48, 1988.
- [Ker98] Norbert L Kerr. Harking: Hypothesizing after the results are known. *Personality and social psychology review*, 2(3):196–217, 1998.
- [KGF19] Maryam Kouchaki, Francesca Gino, and Yuval Feldman. The ethical perils of personal, communal relations: A language perspective. *Psychological Science*, 30(12):1745–1766, 2019.
- [KKdL98] Ita GG Kreft, Ita Kreft, and Jan de Leeuw. *Introducing multilevel modeling*. Sage, 1998.
- [KKH19] Alex Kale, Matthew Kay, and Jessica Hullman. Decision-making under uncertainty in research synthesis: Designing for the garden of forking paths. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2019.
- [KLHO19] Janin Koch, Andrés Lucero, Lena Hegemann, and Antti Oulasvirta. May ai? design ideation with cooperative contextual bandits. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2019.
- [KNH16] Matthew Kay, Gregory L Nelson, and Eric B Hekler. Researcher-centered design of statistics: Why bayesian statistics better fit the culture and incentives of hcii. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 4521–4532. ACM, 2016.
- [Kor19] Joos Korstanje. "anova's three types of estimating sums of squares: don't make the wrong choice!". *Towards Data Science, Medium*, 2019.
- [KPHH11] Sean Kandel, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. Wrangler: Interactive visual specification of data transformation scripts. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 3363–3372. ACM, 2011.

- [KPHH12] Sean Kandel, Andreas Paepcke, Joseph M Hellerstein, and Jeffrey Heer. Enterprise data analysis and visualization: An interview study. *IEEE Transactions on Visualization and Computer Graphics*, 18(12):2917–2926, 2012.
- [KPRP07] Gary Klein, Jennifer K Phillips, Erica L Rall, and Deborah A Peluso. A data-frame theory of sensemaking. In *Expertise out of context*, pages 118–160. Psychology Press, 2007.
- [KR12] Maurits Kaptein and Judy Robertson. Rethinking statistical analysis methods for chi. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1105–1114. ACM, 2012.
- [KS99] David Klahr and Herbert A Simon. Studies of scientific discovery: Complementary approaches and convergent findings. *Psychological Bulletin*, 125(5):524, 1999.
- [KVR20] Max Kuhn, Davis Vaughan, and RStudio. *parsnip: A common api to modeling and analysis functions*, 2020.
- [KW19] Scott Klemmer and Jacob Wobbrock. Designing, running, and analyzing experiments, 2019.
- [KW20] Max Kuhn and Hadley Wickham. *Tidymodels: a collection of packages for modeling and machine learning using tidyverse principles.*, 2020.
- [LA15] Steson Lo and Sally Andrews. To transform or not to transform: Using generalized linear mixed models to analyse reaction time data. *Frontiers in psychology*, 6:1171, 2015.
- [LAH19] Yang Liu, Tim Althoff, and Jeffrey Heer. Paths explored, paths omitted, paths obscured: Decision points & selective reporting in end-to-end data analysis. *arXiv preprint arXiv:1910.13602*, 2019.
- [LAvA⁺19] Minha Lee, Sander Ackermans, Nena van As, Hanwen Chang, Enzo Lucas, and Wijnand IJsselsteijn. Caring for vincent: A chatbot for self-compassion. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- [LBE19] Jiali Liu, Nadia Boukhelifa, and James R Eagan. Understanding the role of alternatives in data analysis practices. *IEEE transactions on visualization and computer graphics*, 26(1):66–76, 2019.
- [LCD⁺19a] Michael D Lee, Amy H Criss, Berna Devezter, Christopher Donkin, Alexander Etz, Fábio P Leite, Dora Matzke, Jeffrey N Rouder, Jennifer S Trueblood, Corey N White, et al. Robust modeling in cognitive science. *Computational Brain & Behavior*, 2(3):141–153, 2019.
- [LCD⁺19b] Merrin Man Long Leong, Arthur Kwok Leung Cheung, Wei Dai, Sai Wah Tsao, Chi Man Tsang, Christopher W. Dawson, Josephine Mun Yee Ko, and Maria Li Lung. Ebv infection is associated with histone bivalent switch modifications in squamous epithelial cells. *116(28):14144–14153*, 2019.

- [LDG⁺14] James Lloyd, David Duvenaud, Roger Grosse, Joshua Tenenbaum, and Zoubin Ghahramani. Automatic construction and natural-language description of non-parametric regression models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014.
- [LGL⁺19] Chang Liu, Hui-Min Ge, Bai-Hui Liu, Rui Dong, Kun Shan, Xue Chen, Mu-Di Yao, Xiu-Miao Li, Jin Yao, Rong-Mei Zhou, Shu-Jie Zhang, Qin Jiang, Chen Zhao, and Biao Yan. Targeting pericyte–endothelial cell crosstalk by circular rna-cpwwp2a inhibition aggravates diabetes-induced microvascular dysfunction. 116(15):7455–7464, 2019.
- [Lib19] Kent State University Libraries. Spss tutorials: Analyzing data, 2019.
- [LLC20a] StataCorp LLC. Language syntax, 2020.
- [LLC20b] StataCorp LLC. Stata 16 documentation, 2020.
- [LP20] Erin LeDell and Sebastien Poirier. H2O AutoML: Scalable automatic machine learning. *7th ICML Workshop on Automated Machine Learning (AutoML)*, July 2020.
- [LTBS00] David J. Lunn, Andrew Thomas, Nicky Best, and David Spiegelhalter. Winbugs - a bayesian modelling framework: Concepts, structure, and extensibility. *Statistics and Computing*, 10(4):325–337, Oct 2000.
- [LTE16] Calvin Loncaric, Emina Torlak, and Michael D Ernst. Fast synthesis of fast collections. *ACM SIGPLAN Notices*, 51(6):355–368, 2016.
- [MABL⁺07] Wendy E Mackay, Caroline Appert, Michel Beaudouin-Lafon, Olivier Chapuis, Yangzhou Du, Jean-Daniel Fekete, and Yves Guiard. Touchstone: exploratory design of experiments. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 1425–1434. ACM, 2007.
- [McE20] Richard McElreath. *Statistical rethinking: A Bayesian course with examples in R and Stan*. CRC press, 2020.
- [McN15] Amelia Ahlers McNamara. *Bridging the gap between tools for learning and for doing statistics*. PhD thesis, UCLA, 2015.
- [McN18] Amelia McNamara. Key attributes of a modern statistical computing tool. *The American Statistician*, 2018.
- [MSF19] Nora McDonald, Sarita Schoenebeck, and Andrea Forte. Reliability and inter-rater reliability in qualitative research: Norms and guidelines for cscw and hci practice. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–23, 2019.
- [MSN⁺20] Arni Magnusson, Hans Skaug, Anders Nielsen, Casper Berg, Kasper Kristensen, Martin Maechler, Koen van Bentham, Ben Bolker, Nafis Sadat, Daniel Lüdecke, Russ Lenth, Joseph O’Brien, and Mollie Brooks. Package ‘glmmtnb’. 2020.

- [MTB⁺19] Philippe Mauffrey, Nicolas Tchitchev, Vilma Barroca, Alexis Bemelmans, Virginie Firlej, Yves Allory, Paul-Henri Romeo, and Claire Magnon. Progenitors from the central nervous system drive neurogenesis in cancer. *Nature*, 569(7758):672–678, 2019.
- [MTM⁺19] Kevin McCloskey, Ankur Taly, Federico Monti, Michael P. Brenner, and Lucy J. Colwell. Using attribution to decode binding mechanism in neural network models for chemistry. 116(24):11624–11629, 2019.
- [Mut23] Muthén & Muthén. Mplus, 2023.
- [MWN⁺19] Dominik Moritz, Chenglong Wang, Greg L Nelson, Halden Lin, Adam M Smith, Bill Howe, and Jeffrey Heer. Formalizing visualization design knowledge as constraints: Actionable and extensible models in draco. *IEEE transactions on visualization and computer graphics*, 25(1):438–448, 2019.
- [Nav19] Danielle Navarro. xprmntr, 2019.
- [Nav21] Danielle Navarro. Jaysire: Building jspysch experiments in r, 2021.
- [NHNO19] Chi T. Ngo, Aidan J. Horner, Nora S. Newcombe, and Ingrid R. Olson. Development of holistic episodic recollection. *Psychological Science*, 30(12):1696–1706, 2019.
- [Nor86] Donald A Norman. Cognitive engineering. *User centered system design*, 31:61, 1986.
- [Nor13] Don Norman. *The design of everyday things: Revised and expanded edition*. Basic books, 2013.
- [NS⁺72] Allen Newell, Herbert Alexander Simon, et al. *Human problem solving*, volume 104. Prentice-Hall Englewood Cliffs, NJ, 1972.
- [NSVW19] Taylor D Nadauld, Berk A Sensoy, Keith Vorkink, and Michael S Weisbach. The liquidity cost of private equity investments: Evidence from secondary market transactions. *Journal of Financial Economics*, 132(3):158–181, 2019.
- [NW72] John Ashworth Nelder and Robert WM Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society: Series A (General)*, 135(3):370–384, 1972.
- [oA20] University of Amsterdam. Jasp: A fresh way to do statistics, 2020.
- [OKB⁺19] Mikhail M Otkrovskikh, Ilya I Klimovskikh, Hendrik Bentmann, D Estyunin, Alexander Zeugner, Ziya S Aliev, S Gaß, AUB Wolter, AV Koroleva, Alexander M Shikin, et al. Prediction and observation of an antiferromagnetic topological insulator. *Nature*, 576(7787):416–422, 2019.
- [Oli06] Travis E Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.
- [ONK⁺19] Luigi Ombrato, Emma Nolan, Ivana Kurelac, Antranik Mavousian, Victoria Louise Bridgeman, Ivonne Heinze, Probir Chakravarty, Stuart Horswell, Estela Gonzalez-Gualda, Giulia Matachione, et al. Metastatic-niche labelling reveals parenchymal cells with stem features. *Nature*, 572(7771):603–608, 2019.

- [P⁺00] Judea Pearl et al. Models, reasoning and inference. *Cambridge, UK: Cambridge University Press*, 19, 2000.
- [PAH⁺19] Emil Persson, Erkin Asutay, Markus Heilig, Andreas Löfberg, Nancy Pedersen, Daniel Västfjäll, and Gustav Tinghög. Variation in the μ -opioid receptor gene (oprm1) does not moderate social-rejection sensitivity in humans. *Psychological Science*, 30(7):1050–1062, 2019.
- [PBD⁺20] José Pinheiro, Douglas Bates, Saikat DebRoy, Deepayan Sarkar, EISPACK authors, Siem Heisterkamp, Bert Van Willigen, and R-core. Package ‘nlme’. 2020.
- [PC05] Peter Pirolli and Stuart Card. The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In *Proceedings of international conference on intelligence analysis*, volume 5, pages 2–4. McLean, VA, USA, 2005.
- [Pea95a] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.
- [Pea95b] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.
- [Pei07] Jonathan W Peirce. Psychopy—psychophysics software in python. *Journal of neuroscience methods*, 162(1-2):8–13, 2007.
- [Pfa97] M Pfannkuch. Statistical thinking: One statistician’s perspective. *Research papers on stochastics education*, pages 171–178, 1997.
- [PGE⁺18] Pavel Panchevka, Adam T Geller, Michael D Ernst, Zachary Tatlock, and Shoaib Kamil. Verifying that web pages have accessible layout. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 1–14. ACM, 2018.
- [PGEE⁺19] Anthony Poon, Sarah Giroux, Parfait Eloundou-Enyegue, François Guimbretière, and Nicola Dell. Engaging high school students in cameroon with exam practice quizzes via sms and whatsapp. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- [PGSF19] Liuba Papeo, Nicolas Goupil, and Salvador Soto-Faraco. Visual search for people among people. *Psychological Science*, 30(10):1483–1496, 2019.
- [Pra19] Tanu N Prabhu. Top python libraries used in data science. *Towards Data Science, Medium*, 2019.
- [PSTsd20] Josef Perktold, Skipper Seabold, Jonathan Taylor, and statsmodels developers. Statsmodels v0.10.2 reference guide. 2020.
- [PVG⁺11] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.

- [PW⁺00] Maxine Pfannkuch, Chris J Wild, et al. Statistical thinking an statistical practice: Themes gleaned from professional statisticians. *Statistical science*, 15(2):132–152, 2000.
- [RSPC93] Daniel M Russell, Mark J Stefk, Peter Pirolli, and Stuart K Card. The cost structure of sensemaking. In *Proceedings of the INTERACT’93 and CHI’93 conference on Human factors in computing systems*, pages 269–276. ACM, 1993.
- [Rub04] Donald B Rubin. Teaching statistical inference for causal effects in experiments and observational studies. *Journal of Educational and Behavioral Statistics*, 29(3):343–367, 2004.
- [RVB⁺20] Brian Ripley, Bill Venables, Douglas M. Bates, Kurt Hornik, Albrecht Gebhardt, , and David Firth. Package ‘mass’. 2020.
- [SAS20a] SAS. Jmp, 2020.
- [SAS20b] SAS. Jmp, 2020.
- [SG89] Susan Leigh Star and James R Griesemer. Institutional ecology,translations’ and boundary objects: Amateurs and professionals in berkeley’s museum of vertebrate zoology, 1907-39. *Social studies of science*, 19(3):387–420, 1989.
- [SH17] N.J.A. Sloane and R.H. Hardin. Gosset: A general-purpose program for designing experiments, 2017.
- [Sha10] William R Shadish. Campbell and rubin: A primer and comparison of their approaches to causal inference in field settings. *Psychological methods*, 15(1):3, 2010.
- [SK95] Christian D Schunn and David Klahr. A 4-space model of scientific discovery. In *Proceedings of the 17th annual conference of the cognitive science society*, pages 106–111, 1995.
- [SK96] Christian D Schunn and David Klahr. When and how to go beyond a 2-space model of scientific discovery. In *Proceedings of the Eighteenth Annual Conference of the Cognitive Science Society: July 12–15, 1996, University of California, San Diego*, volume 18, page 25. Psychology Press, 1996.
- [SKF18] Amanda Sweeny, Andrew J Ko, and James Fogarty. Scout: Mixed-initiative exploration of design variations through high-level design constraints. In *The 31st Annual ACM Symposium on User Interface Software and Technology Adjunct Proceedings*, pages 134–136. ACM, 2018.
- [sld20] scikit-learn developers. Scikit-learn v0.23.2 documentation, 2020.
- [Smo19] Michael Smolyansky. Policy externalities and banking integration. *Journal of Financial Economics*, 132(3):118–139, 2019.
- [SMWH17] Arvind Satyanarayanan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. Vega-lite: A grammar of interactive graphics. *IEEE transactions on visualization and computer graphics*, 23(1):341–350, 2017.

- [SOR⁺19a] Joon-Gi Shin, Eiji Onchi, Maria Jose Reyes, Junbong Song, Uichin Lee, Seung-Hee Lee, and Daniel Saakes. Slow robots for unobtrusive posture correction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–10, 2019.
- [SOR⁺19b] Joon-Gi Shin, Eiji Onchi, Maria Jose Reyes, Junbong Song, Uichin Lee, Seung-Hee Lee, and Daniel Saakes. Slow robots for unobtrusive posture correction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–10, 2019.
- [SP10] Skipper Seabold and Josef Perktold. Statsmodels: Econometric and statistical modeling with python. In *Proceedings of the 9th Python in Science Conference*, volume 57, page 61. Scipy, 2010.
- [Spi94] Peter Spirtes. Conditional independence in directed cyclic graphical models for feedback. 1994.
- [SPS21] IBM SPSS. Spss software, 2021.
- [SRM⁺96] Peter Spirtes, Thomas Richardson, Christopher Meek, Richard Scheines, and Clark Glymour. Using d-separation to calculate zero partial correlations in linear models with correlated errors. *Publisher: Carnegie Mellon University*, 1996.
- [SSY20] Etsushi Suzuki, Tomohiro Shinozaki, and Eiji Yamamoto. Causal diagrams: pitfalls and tips. *Journal of epidemiology*, page JE20190192, 2020.
- [Sta21] Stata. Stata software, 2021.
- [Suh14] Michael Suh. Higher education, gender & work dataset, 2014.
- [SV18] Etsushi Suzuki and Tyler J VanderWeele. Mechanisms and uncertainty in randomized controlled trials: A commentary on deaton and cartwright. *Social science & medicine (1982)*, 210:83–85, 2018.
- [SWF16] John Salvatier, Thomas V Wiecki, and Christopher Fonnesbeck. Probabilistic programming in python using pymc3. *PeerJ Computer Science*, 2:e55, 2016.
- [Tan21] Emi Tanaka. Edibble: An r-package to construct designs using the grammar of experimental design, 2021.
- [Tcw20] R Core Team and contributors worldwide. Package ‘stats’ v4.1.0. *CRAN*, 2020.
- [THHLB13] C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown. Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms. In *Proc. of KDD-2013*, pages 847–855, 2013.
- [THK11] Johannes Textor, Juliane Hardt, and Sven Knüppel. Dagitty: a graphical tool for analyzing causal diagrams. *Epidemiology*, 22(5):745, 2011.
- [TM20a] Inc. The MathWorks. Matlab, 2020.

- [TM20b] Inc. The MathWorks. Statistics and machine learning toolbox. 2020.
- [TPO⁺19] Nina Thigpen, Nathan M. Petro, Jessica Oschwald, Klaus Oberauer, and Andreas Keil. Selection of visual objects in perception and working memory one at a time. *Psychological Science*, 30(9):1259–1272, 2019.
- [VAK18] Lea Verou, Tarfah Alrashed, and David Karger. Extending a reactive expression language with data update actions for end-user application authoring. In *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*, pages 379–387, 2018.
- [Van87] Walter Vandaele. *Participation in illegitimate activities: Ehrlich revisited, 1960*, volume 8677. Inter-university Consortium for Political and Social Research, 1987.
- [Van19] Tyler J VanderWeele. Principles of confounder selection. *European journal of epidemiology*, 34(3):211–219, 2019.
- [VD00] András Vargha and Harold D Delaney. A critique and improvement of the cl common language effect size statistics of mcgraw and wong. *Journal of Educational and Behavioral Statistics*, 25(2):101–132, 2000.
- [VDN⁺13] Priscilla Velentgas, Nancy A Dreyer, Parivash Nourjah, Scott R Smith, Marion M Torchia, et al. Developing a protocol for observational comparative effectiveness research: a user’s guide. 2013.
- [VR13] William N Venables and Brian D Ripley. *Modern applied statistics with S-PLUS*. Springer Science & Business Media, 2013.
- [VTB⁺19] Dario Veneziano, Luisa Tomasello, Veronica Balatti, Alexey Palamarchuk, Laura Z. Rassenti, Thomas J. Kipps, Yuri Pekarsky, and Carlo M. Croce. Dysregulation of different classes of trna fragments in chronic lymphocytic leukemia. 116(48):24252–24258, 2019.
- [W⁺14] Hadley Wickham et al. Tidy data. *Journal of Statistical Software*, 59(10):1–23, 2014.
- [WAB⁺19] Hadley Wickham, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Grolemund, Alex Hayes, Lionel Henry, Jim Hester, et al. Welcome to the tidyverse. *Journal of Open Source Software*, 4(43):1686, 2019.
- [Wik19a] Wikipedia contributors. Jmp (statistical software) — Wikipedia, the free encyclopedia. [https://en.wikipedia.org/w/index.php?title=JMP_\(statistical_software\)&oldid=887217350](https://en.wikipedia.org/w/index.php?title=JMP_(statistical_software)&oldid=887217350), 2019. [Online; accessed 5-April-2019].
- [Wik19b] Wikipedia contributors. R (programming language) — Wikipedia, the free encyclopedia. [https://en.wikipedia.org/w/index.php?title=R_\(programming_language\)&oldid=890657071](https://en.wikipedia.org/w/index.php?title=R_(programming_language)&oldid=890657071), 2019. [Online; accessed 5-April-2019].

- [Wik19c] Wikipedia contributors. Sas (software) — Wikipedia, the free encyclopedia. [https://en.wikipedia.org/w/index.php?title=SAS_\(software\)&oldid=890451452](https://en.wikipedia.org/w/index.php?title=SAS_(software)&oldid=890451452), 2019. [Online; accessed 5-April-2019].
- [Wik19d] Wikipedia contributors. Spss — Wikipedia, the free encyclopedia. <https://en.wikipedia.org/w/index.php?title=SPSS&oldid=888470477>, 2019. [Online; accessed 5-April-2019].
- [Wil99] Leland Wilkinson. Statistical methods in psychology journals: Guidelines and explanations. *American psychologist*, 54(8):594, 1999.
- [WLH19] Kanit Wongsuphasawat, Yang Liu, and Jeffrey Heer. Goals, process, and challenges of exploratory data analysis: An interview study. *arXiv preprint arXiv:1911.00568*, 2019.
- [WP99] Chris J Wild and Maxine Pfannkuch. Statistical thinking in empirical enquiry. *International statistical review*, 67(3):223–248, 1999.
- [WSVB15] Chat Wacharamoortham, Krishna Subramanian, Sarah Theres Volkel, and Jan Borchers. Statsplorer: Guiding novices in statistical analysis. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pages 2693–2702. ACM, 2015.
- [WWL21] R Michael Winters, Bruce N Walker, and Grace Leslie. Can you hear my heartbeat?: Hearing an expressive biosignal elicits empathy. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–11, 2021.
- [YK20] Bin Yu and Karl Kumbier. Veridical data science. *Proceedings of the National Academy of Sciences*, 117(8):3920–3929, 2020.
- [YWV⁺19] Chun Yu, Xiaoying Wei, Shubh Vachher, Yue Qin, Chen Liang, Yueling Weng, Yizheng Gu, and Yuanchun Shi. Handsee: Enabling full hand interaction on smartphone with front camera-based stereo vision. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- [ZHC19] Xiaobo Zhao, Jianyan Huang, and Joanne Chory. Gun1 interacts with morf2 to regulate plastid rna editing during retrograde signaling. *Proceedings of the National Academy of Sciences*, 116(20):10162–10167, 2019.
- [ZTDB19] Anqing Zheng, Elliot M. Tucker-Drob, and Daniel A. Briley. National gross domestic product, science interest, and science achievement: A direct replication and extension of the tucker-drob, cheung, and briley (2014) study. *Psychological Science*, 30(5):776–788, 2019.
- [ZWY⁺19] Zixu Zhou, Qiuyang Wu, Zhangming Yan, Haizi Zheng, Chien-Ju Chen, Yuan Liu, Zhijie Qi, Riccardo Calandrelli, Zhen Chen, Shu Chien, H. Irene Su, and Sheng Zhong. Extracellular rna in a single droplet of human serum reflects physiologic and disease states. 116(38):19200–19208, 2019.

Chapter A

Appendix A: Content analysis resources

This appendix provides greater detail describing the corpus of research publications we scraped and analyzed to conduct our content analysis (Chapter 4, Section 4.2), our coding procedure, codebook (Table A.1), and a summary for each paper.

A.1 Dataset overview

To collect our corpus of research publications, we scraped the five venues' proceedings using *Heleena* [CMB18] and wrote additional scripts to randomly sample from each venue. The first author read papers and included them in the sample if they used statistical analyses. We did not discriminate between papers that used statistical analyses as a primary or secondary methodology. We anticipated that authors would describe hypothesis formalization in both cases.

We coded a total of 2,989 paragraphs across 50 papers. Results were the most commonly discussed topic. Approximately 31% of the paragraphs (in 50 papers) discussed interpretations of statistical results, and 11% (in 37 papers) provided details about statistical results (e.g., parameter estimates). Interpreted results often co-occurred with statistical results. 21% of paragraphs (in 40 papers) described data collection design (e.g., how the experiment was designed, how the data were collected, etc.). Specifications of statistical models appeared in 19% of paragraphs (in 50 papers). 11% of paragraphs (in 45 papers) discussed proxy variables, or measures to quantify abstract constructs (e.g., music enjoyment).

Researchers mentioned software used for statistical analysis in 3% of paragraphs (in 25 papers), sometimes even specifying function names and parameters, a level of detail we did not expect to find in publications. To our surprise, more papers mentioned software than included equations. Only fifteen papers (JFE: 9, PS: 5, PNAS: 1) included equations in a total of 71 paragraphs. This suggests that mathematical equations, though part of the hypothesis formalization process, are less important to researchers than their tool-specific implementations.

A.2 Procedure

Based on exploratory rounds of open coding on an independent sample and noticeable differences in writing structure and style across the venues, we read all sections of papers unlike McDonald et al. [MSF19] who performed a similar content analysis of qualitative analysis methodologies and focused on methods sections only. We also read the materials and methods section(s) included after

references in papers from the PNAS and Nature venues, but otherwise we did not code any figures, tables, and auxiliary materials.

The first and second author developed the codebook, analyzed five papers (one from each venue), discussed agreements and disagreements, and iterated on the analysis protocol and codebook. The first two authors then used the revised code book to analyze another two papers that were substantially different in data analysis approach and writing style, discussed any disagreements, and refined the codebook. The coders reached substantial agreement ($IRR = .69 - .72$) even before resolving disagreements.

We coded the papers at the paragraph-level. We initially started by coding at the sentence-level but found that paragraphs provided necessary context for accurately interpreting and coding sentences, showed co-occurrence patterns, and were more expedient and anecdotally more reliable to code. Nonetheless, throughout the coding process, we deliberated and discussed key sentences in paragraphs that shaped the paper’s argumentation structure. The codebook contains such key sentences.

After the first three authors coded, reviewed (for coding consistency), and discussed each paper, we created “reorderable matrices” [Ber11] for each paper. The first three authors scrutinized the matrices and cross-referenced the matrices and papers to identify a set of visual patterns. The visual patterns indicated how researchers structured their scientific arguments (e.g., Pattern 1); specified and summarized research questions and hypotheses, indicative of hypothesis refinement (e.g., Pattern 2); decomposed their research questions and hypotheses from more general to more specific ones, indicative of hypothesis refinement (e.g., Pattern 3); described their data collection and cleaning procedures, sometimes also discussing specific proxies (e.g., Patterns 4, 5); mentioned software and computational settings relevant to model implementation details (e.g., Pattern 6); and discussed statistical specifications and results (e.g., Patterns 7). The definitions and notes on the patterns we used are included as supplementary material. Please refer to the README for file names and descriptions.

A.3 Codebook

A.4 Additional findings: Contribution types

We identified papers that presented empirical findings (41 papers), validated a prototype system (8 papers), or developed a new methodology (6 papers).

After reading and coding the papers, we re-read assigned each paper at least one of the following contribution types: Methodology, System or Technique, and Empirical Findings, and Other. Methodological contributions introduce a new way of measuring a concept and may be in the form of novel experimental designs, procedures, proxies, or other measures. System or technique contributions develop a prototype tool, which may be physical, biological, or chemical in nature. Empirical findings contributions primarily show or explain a new phenomenon, which may involve developing new causal models of a domain. Other contributions included replication studies and other results that were unique to one or two papers in our sample, such as finding a new species in [DMC⁺19]. We identified these four contribution types through discussions and open coding.

We found that 41 papers that made empirical contributions describing or explaining a phenomenon; eight papers that developed and evaluated physical or biological prototype tools; and six papers that presented novel methodologies such as experimental protocols or measures. Ten

Table A.1: The codebook for analyzing the content of research publications.

| Codes | Definitions and Examples | %Occurrences |
|-------|--------------------------|--------------|
|-------|--------------------------|--------------|

| | | |
|---|---|------|
| Research Goals | | |
| Question or statement of unknown Predicted outcomes | Explicit, clear statement about an unknown phenomenon or an open-ended question “However, the ontogeny of holistic recollection is uncharted.” [NHNO19] | 4.2 |
| Specific statistical expectation | A clear conjecture of an outcome that does not specify a specific mathematical relationship “We hypothesized that the outward current is mainly carried by FA anions...” [BCK ⁺ 19] | 6.8 |
| Specific objectives | A conjecture specifying how observations will be related to one another mathematically/statistically “If this dependency measure (data-independent model) was significantly greater than zero, this provided evidence for significant retrieval dependency...” [NHNO19] | 0.1 |
| Examination of associations | Statements about reaching objectives “To assess the potential clinical relevance of the neo-development of a neuronal network in prostate cancer, DCX+ cells were quantified in benign prostate hyperplasia...” [MTB ⁺ 19] | 2.6 |
| Data Sample Information | | |
| Study design and protocol | Information about the procedures used or prototypes developed to collect the data for analysis, such as any assays or experimental designs, including any limitations (e.g., conditions/randomization, interventions, treatments) “Before the experiment, we introduced the working principle of HandSee. Then we tested the techniques one by one. For each technique, we first demonstrated our interaction technique. After...” [Y WV ⁺ 19] | 20.8 |
| Initial data sourcing | Information about the source, size, and characteristics of the data sample that was collected or analyzed “A total of 32 four-year-old children (15 female; age: $M = 52.05$ months, $SD = 3.37$) and 30 six-year-old children (17 female; age: $M = 76.37$ months; $SD = 2.16$) from the Philadelphia area participated in the study...” [NHNO19] | 9.8 |
| Data filtering/sampling | Any criteria, procedures, and decisions to filter, remove, combine, and split data for data quality, sub-analyses, or robustness (e.g., sampling from existing datasets, removing outliers, etc.) “As for the US data, we restrict our attention to sectors with ten or more firms. ” [GSY19] | 6.4 |
| Details about data used for analysis | Any summary statistics (e.g., mean, standard deviation, distribution, etc.) and other information describing the final data sample used for statistical analysis “All subjects’ mean values were within 2.5 standard deviations from the group mean; therefore, they were all included in the following analyses.” [PGSF19] | 2.9 |
| Statistical Analysis | | |
| Proxy | Any information about how concepts are measured, including any limitations, etc.; Can be established or new ways of measuring a construct “Our definition of a price run-up is based on the industry value-weighted return.” [GSY19] | 12.3 |
| Equation | Any mathematical equation, using symbols or sentences “The absolute number of cells was calculated as ((number of Lin-eYFP+ cells acquired cellularity of the organ)/number of live single cells acquired).” [MTB ⁺ 19] | 2.4 |
| Statistical specification | Describing a statistical model (e.g., linear regression), test (e.g. Student’s t-test), or other procedure (e.g., contingency table) for analyzing the data “Frequentist null-hypothesis significance testing was complemented with Bayesian hypothesis testing, which quantified the evidence for the presence or absence of effects...” [BLT ⁺ 19] | 18.7 |
| Results | | |
| Statistical results | Reporting the findings (not usage) of statistical analyses or models, refer to specific quantified metrics (e.g., ratio, coefficient, correlation, etc.) with specific values (e.g., numbers) or aspects of values (e.g., positive estimate, positive relationship) “In all experiments, when the entire sample size ($N = 24$) was included in the analyses, the main findings in each experiment remained significant for all color-memory estimates (for paired comparisons, all t s > 12.74, p s .012, and $BF10$ s = 4.32; for three-group comparisons, all F s > 7.07, p s .0021, and $BF10$ s = 16.64). ” [BLT ⁺ 19] | 10.7 |
| Interpreted results | What the statistical results mean conceptually “This result supports the notion that the economies of scale...can induce larger firms to hedge more extensively.” [HHZ19] | 30.9 |
| Causal model | A causal model or mechanism (with a clear cause and effect) supported by the data and statistical analysis results “Here we show that CUN1 interacts with MORE0/RIPD, hence in only the name MORE0 | 0.2 |

papers made various other contributions (e.g., replicating a previous study, finding a new species, developing a design space, etc.). Tables A.2 through A.6 give an overview of contribution types in each venue. We separated the tables by venue due to spacing constraints.

Papers contributing empirical findings consisted of ten papers from PNAS, ten from PS, eight from JFE, eight from Nature, and five from CHI. Six of the eight system/technique contributions came from CHI papers, with one each from Nature and PNAS. Out of the six methodology contributions, three came from JFE papers, two from CHI, and one from PNAS. Thirteen papers fell under multiple contribution types. Co-occurrences of two out of the three contribution types were seen in a few of the CHI and PNAS papers, with system/technique contributions co-occurring with either methodology or empirical findings. Co-occurrences in the PS and PNAS papers involved an "Other" contribution type occurring most often with empirical findings. We identified only one JFE paper with multiple contributions; in this case, methodology and empirical findings co-occurred. We did not notice any obvious differences in paper content or structure due to research contribution types, either within or across venues.

A.5 Summary of papers analyzed

Table A.2: Summary of CHI papers in our dataset.

| Title | Short summary | Method | System | Empirical | Other |
|--|---|--------|--------|-----------|-------|
| ACM Conference on Human Factors in Computing Systems (CHI). | | | | | |
| Detecting Visuo-Haptic Mismatches in Virtual Reality using the Prediction Error Negativity of Event-Related Brain Potentials [GAL ⁺ 19] | The authors develop a new, more objective metric for haptic immersion. Through a user study, they find that the new metric is able to detect visuo-haptic mismatches in VR. | ✓ | — | — | — |
| Engaging High School Students in Cameroon with Exam Practice Quizzes via SMS and WhatsApp [PGEE ⁺ 19] | The researchers provide study support through quiz questions delivered through SMS or WhatsApp. The researchers observe differences in participation during a three-week deployment. | — | — | ✓ | — |
| Springlets: Expressive, Flexible and Silent On-Skin Tactile Interfaces [HWV ⁺ 19] | Springlets is a mechano-tactile interface for skin. The authors discuss its design, fabrication, user perceptions, and possible applications. | — | ✓ | — | — |
| HandSee: Enabling Full Hand Interaction on Smartphones with Front Camera-based Stereo Vision [YWV ⁺ 19] | The authors are able to detect a phone user's gripping and touching interactions by using a mirror on the front camera to obtain stereo vision. | ✓ | ✓ | — | ✓ |
| User Attitudes towards Algorithmic Opacity and Transparency in Online Reviewing Platforms [EVL ⁺ 19] | The authors ask three research questions around how Yelp users view algorithmic control/intervention on the platform and how increasing their awareness of it through a system changes their perspectives and opinions about it. The authors find that individuals that are more invested in Yelp as reviewers are more likely to defend the platform's algorithms. | — | — | ✓ | — |
| Slow Robots for Unobtrusive Posture Correction [SOR ⁺ 19a] | The authors develop and evaluate a system for automatically correcting user posture. The authors conduct two empirical studies, one formative and one evaluation. The formative study identifies end user perception of moving screens. The evaluation study evaluates user experience and how often users corrected their posture. The authors use mixed methods, both quantitative and qualitative/observational in their studies. | — | ✓ | — | — |
| May AI? Design Ideation with Cooperative Contextual Bandits [KLHO19] | The authors develop a new technique and system for co-creation with an AI. They evaluate the effects of the AI on a series of task and creativity measures. They find that the AI is helpful in some dimensions. | — | ✓ | — | — |
| The Effects of Interruption Timings on Autonomous Height-Adjustable Desks that Respond to Task Changes [SOR ⁺ 19b] | The authors investigate the most opportune time to adjust desks for improved ergonomics while minimizing interruption and annoyance/negative experiences during tasks. The authors find that changing desk height during task transition periods are the least disruptive, but end-users are dubious/have less trust in the automated adjustments. On the other hand, adjusting desk height after end-users have initiated a new task/changed tasks causes increased disruption but also increased trust in the automated desk. | — | — | ✓ | — |
| Caring for Vincent: A Chatbot for Self-compassion [LAvA ⁺ 19] | The authors design a chatbot and then see how taking care of or being taken by the chatbot affects self-compassion. The authors test this hypothesis quantitatively and then follow-up with additional analyses about gender and tendency for self-compassion. The authors further contextualize these quantitative results with qual- | — | ✓ | ✓ | — |

Table A.3: Summary of JFE papers in our dataset.

| Title | Short summary | Method | System | Empirical | Other |
|--|---|--------|--------|-----------|-------|
| Journal of Financial Economics (JFE) | | | | | |
| The relevance of broker networks for information diffusion in the stock market [DMFKS19] | The authors find evidence for the spread of information from central brokers to their best clients/investors to more peripheral clients/investors, which benefit the central brokers through high returns. There are three main findings: (i) more central brokers have higher returns, (ii) this can be seen/validated in how informed investors trade, and (iii) information affects “price discovery.” | — | — | ✓ | — |
| Do firms hedge with foreign currency derivatives for employees? [HZH19] | The authors present evidence of a relationship between firms’ employee treatment scores and the fraction of revenue hedged with currency-based derivatives. | — | — | ✓ | — |
| The liquidity cost of private equity investments: Evidence from secondary market transactions [NSVW19] | The authors characterize how transaction costs in the secondary market for private equity stakes are determined. | — | — | ✓ | — |
| Institutional investor cliques and governance [CKM19] | The authors examine how investors coordinate to influence governance. | ✓ | — | ✓ | — |
| Policy externalities and banking integration [Smo19] | The author investigates if and how policies aimed at the banking sector in one region have ripple effects in other regions. The author identifies “financial linkages” as the “transmission channel” for these policies to affect other areas. | — | — | ✓ | — |
| Do labor markets discipline? Evidence from RMBS bankers [GKM19] | The authors examine if and how there were any disciplinary actions taken in the mortgage/housing banking industry after the housing crisis of 2008. The authors consider internal and external discipline (within firms and overall). They find that there were no disciplinary measures taken. | — | — | ✓ | — |
| Firing the wrong workers: Financing constraints and labor misallocation [CCM19a] | The authors develop a theoretical model for the impact of financial constraints on firing. They find that firms first fire short-term workers even though those workers might provide longer-term value to the firms. | ✓ | — | — | — |
| Time-varying ambiguity, credit spreads, and the levered equity premium [CCM19b] | The authors propose a new proxy of ambiguity that captures Knightian uncertainty. They build a new model using this proxy to explain credit spreads and pricing of equity and other financial metrics. | ✓ | — | — | — |
| The power of shareholder votes: Evidence from uncontested director elections [ADP19] | The authors investigate the impact of dissenting votes on directors who are elected without contestation. The authors find that shareholder votes impact director’s career trajectories negatively. | — | — | ✓ | — |
| Bubbles for Fama [GSY19] | The authors test a widely held theory/hypothesis about if stock prices experience bubbles and can be detected a priori. They find that they cannot (cannot disprove hypothesis) and find that sharp price increases predict a probability of a crash, and several other factors predict future crashes and returns. | — | — | ✓ | — |

Table A.4: Summary of Nature papers in our dataset.

| Title | Short summary | Method | System | Empirical | Other |
|---|--|--------|--------|-----------|-------|
| Nature | | | | | |
| VISTA is an acidic pH-selective ligand for PSGL-1 [JSP ⁺ 19] | The authors find that V-domain immunoglobulin suppressor of T cell activation (VISTA) suppress T cells in acidic pH environments, including tumor microenvironments. | — | — | ✓ | — |
| A new species of Homo from the Late Pleistocene of the Philippines [DMC ⁺ 19] | The authors discover and analyze bones that they conclude to be a new species, which they call <i>Homo luzonensis</i> . | — | — | ✓ | — |
| H ⁺ transport is an integral function of the mitochondrial ADP/ATP carrier [BCK ⁺ 19] | The authors discover two transport modes for ADP/ATP in mitochondria that explain how energy conversion occurs in mitochondria. | — | — | ✓ | — |
| Antarctic offshore polynyas linked to Southern Hemisphere climate anomalies [CWM ⁺ 19] | Researchers find that polynyas, “large openings in the winter sea ice cover,” develop because of simultaneous upper-ocean preconditioning and meteorological changes. They predict that global warming will continue to create conditions under which polynyas occur. | — | — | ✓ | — |
| Metastatic-niche labelling reveals parenchymal cells with stem features [ONK ⁺ 19] | The authors present a system where metastatic cancer cells “stain” surrounding tissue cells so that researchers can learn about the local cancer environment. The system may enable new discoveries. | — | ✓ | — | ✓ |
| Multi-omics profiling of mouse gastrulation at single-cell resolution [ACM ⁺ 19] | The authors discover the process by which three germ layers develop and differentiate during gastrulation. | — | — | ✓ | — |
| Dynamics and genomic landscape of CD8 ⁺ T cells undergoing hepatic priming [BDSDL ⁺ 19] | The authors identify a cellular reproduction mechanism leveraging the liver (novel) that seems to boost the immune system reactions among HBV patients. | — | — | ✓ | ✓ |
| Prediction and observation of an antiferromagnetic topological insulator [OKB ⁺ 19] | The authors develop a theory about a compound based on measurements of its properties and then use follow-up experiments to test the theory. The authors employ simulations to develop their theory and then a series of experiments that triangulate and test the theorized properties. | — | — | ✓ | — |
| Insect egg size and shape evolve with ecology but not developmental rate [CDdME19] | The authors test three (main) hypotheses in the literature about the factors influencing egg size among insects. The authors find that ecology (where an egg is laid) predicts egg size rather than previously believed-in universal constraints. | — | — | ✓ | ✓ |
| Progenitors from the central nervous system drive neurogenesis in cancer [MTB ⁺ 19] | The authors identify the role of nerves in cancer cell neurogenesis and develop a new model of cancer cell neurogenesis in prostate tumors. Their model incorporates “crosstalk” between the central nervous system and the prostate tumors. They develop this model through the identification of associations between cell groups (“lower-level”) and in mice/humans (host). Their model challenges existing models of cancer. | — | — | ✓ | — |

Table A.5: Summary of PNAS papers in our dataset.

| Title | Short summary | Method | System | Empirical | Other |
|--|--|--------|--------|-----------|-------|
| Proceedings of the National Academy of Sciences (PNAS) | | | | | |
| GUN1 interacts with MORF2 to regulate plastid RNA editing during retrograde signaling [ZHC19] | The authors find that GUN1 and MORF2 affect retrograde signaling and plastid RNA-editing in chloroplasts in plant cells. These findings also suggest that retrograde signaling and plastid RNA editing may be related processes. | — | — | ✓ | — |
| Brain-wide genetic mapping identifies the indusium griseum as a prenatal target of pharmacologically unrelated psychostimulants [FRG ⁺ 19] | The authors find the effects of psychostimulants on fetal development. They find exposure can delay specific kinds of cellular and regional development that may impact child behavior. | — | — | ✓ | — |
| Dysregulation of different classes of tRNA fragments in chronic lymphocytic leukemia [VTB ⁺ 19] | The authors find how two different classes of RNAs are associated with CLL, a type of leukemia most prevalent among humans. Based on their findings, the authors conclude that these classes of RNAs may influence the development of CLL. | — | — | ✓ | — |
| A critical role for microglia in maintaining vascular integrity in the hypoxic spinal cord [HM19] | Through a series of experiments, the authors identify the response/roles of microglia in maintaining the health of a hypoxic spinal cord. The findings suggest microglia's importance in Central Nervous System vascular health. | — | — | ✓ | — |
| SDS22 selectively recognizes and traps metal-deficient inactive PP1 [CMR ⁺ 19] | The authors investigate how SDS22 can both inhibit and activate PP1 (an enzyme). The authors identify a mechanism for SDS22 that explains its behavior. | — | — | ✓ | — |
| Monitoring of switches in heterochromatin-induced silencing shows incomplete establishment and developmental instabilities [BHM19] | The authors were interested in investigating what determines/explains the Position Effect Variegation (PEV) in heterochromatin. Through both mathematical modeling and empirical studies, the researchers find that gene silencing that influences PEV occurs early in embryogenesis but is not stable and changes throughout development. | — | — | ✓ | — |
| Extracellular RNA in a single droplet of human serum reflects physiologic and disease states [ZWY ⁺ 19] | The authors develop and test a new method for sequencing RNAs directly on cell serums using “complementary DNA (cDNA).” The authors test the method/tool’s validity by examining how it can differentiate among many different characteristics in the data—sex, cancer, etc. | — | ✓ | ✓ | — |
| EBV infection is associated with histone bivalent switch modifications in squamous epithelial cells [LCD ⁺ 19b] | Epstein-Barr virus (EBV) infection occurs with some cancers. The authors find evidence that suggests that EBV infection may be related to changes in epithelial cells. | — | — | ✓ | — |
| Targeting pericyte-endothelial cell crosstalk by circular RNA-cPWWP2A inhibition aggravates diabetes-induced microvascular dysfunction [LGL ⁺ 19] | The authors find a mechanism between two different cell “types” that are affected by diabetes. The mechanism suggests new therapeutic interventions for diabetes. | — | — | ✓ | ✓ |
| Using attribution to decode binding mechanism in neural network models for chemistry [MTM ⁺ 19] | The authors develop a metric/process for using “Attribution” as a way to de-bias ML models for learning causal relationships between molecules and binding behaviors. | ✓ | — | ✓ | ✓ |

Table A.6: Summary of PS papers in our dataset.

| Title | Short summary | Method | System | Empirical | Other |
|---|--|--------|--------|-----------|-------|
| Psychological Science | | | | | |
| Working Memory Has Better Fidelity Than Long-Term Memory: The Fidelity Constraint Is Not a General Property of Memory After All [BLT ⁺ 19] | The authors replicate a previous study that found that working memory and long-term memory had identical “fidelity.” The authors find evidence to suggest that this is not the case. | — | — | ✓ | ✓ |
| Separate Contribution of Striatum Volume and Pitch Discrimination to Individual Differences in Music Reward [HAPGNN ⁺ 19] | The authors asked if people’s enjoyment of music is related to neurological structure and ability. The authors find how enjoyment/reward, structure, and ability are related. | — | — | ✓ | — |
| Information Processing Under Reward Versus Under Punishment [BSCN19] | The authors investigate how punishment and reward incentives affect decision making. The authors find that punishment incentives negatively impact decision making. | — | — | ✓ | — |
| Paying Back People Who Harmed Us but Not People Who Helped Us: Direct Negative Reciprocity Precedes Direct Positive Reciprocity in Early Development [CLD ⁺ 19] | The authors asked when and how children learn reciprocity, a key aspect of social coordination. They find that direct negative reciprocity (“paying back” harm) develops earlier than positive reciprocity (“paying back” good), which is generalized rather than directed until children learn social norms. | — | — | ✓ | — |
| Selection of Visual Objects in Perception and Working Memory One at a Time [TPO ⁺ 19] | The authors investigate how things we have seen before (in visual working memory) affect how we perceive what we see now (visual environment). The authors find that humans pay attention to visual aspects that are consistent with memory (“memory-relevant”) and that this processing occurs sequentially in the presence of multiple visual stimuli. | — | — | ✓ | — |
| Variation in the μ -Opioid Receptor Gene (<i>OPRM1</i>) Does Not Moderate Social-Rejection Sensitivity in Humans [PAH ⁺ 19] | Using a much larger sample size and more experimental controls, authors conduct a “conceptual replication” of prior work examining the relationship between a gene and feelings of social-rejection. The authors also provide empirical evidence/test a hypothesis extending prior work. | — | — | ✓ | ✓ |
| The Ethical Perils of Personal, Communal Relations: A Language Perspective [KGF19] | The authors find a link between the warmth of language used and dishonest/cheating behavior. The authors use both controlled experiments and a survey to test mechanisms of this link. | — | — | ✓ | — |
| Visual Search for People Among People [PGSF19] | The authors ask if there is a perceptual unit or mechanism that differentiates between interacting dyads and not interacting dyads. They find evidence for some fundamental perceptual grouping unit that makes grouping interacting/facing dyads easier and individuating interacting dyads harder. | — | — | ✓ | — |
| National Gross Domestic Product, Science Interest, and Science Achievement: A Direct Replication and Extension of the Tucker-Drob, Cheung, and Briley (2014) Study [ZTDB19] | The authors replicate a previous study that found connections among science interest, science achievement, national wealth, and other national characteristics with more recent data. | — | — | ✓ | ✓ |
| Development of Holistic Episodic Recollection [NHNO19] | The authors aim to provide further detail about episodic memory development in humans. Holistic episodic recollection is part of episodic memory, but its development is unknown. They find that holistic recollection increased from 4 to adulthood, finding that 6-year-olds exhibit memory retrieval that is similar to adults despite being less | — | — | ✓ | — |

Chapter B

Appendix B: Tisane's First Release: Additional Examples

B.1 Additional examples of graphs that may be constructed

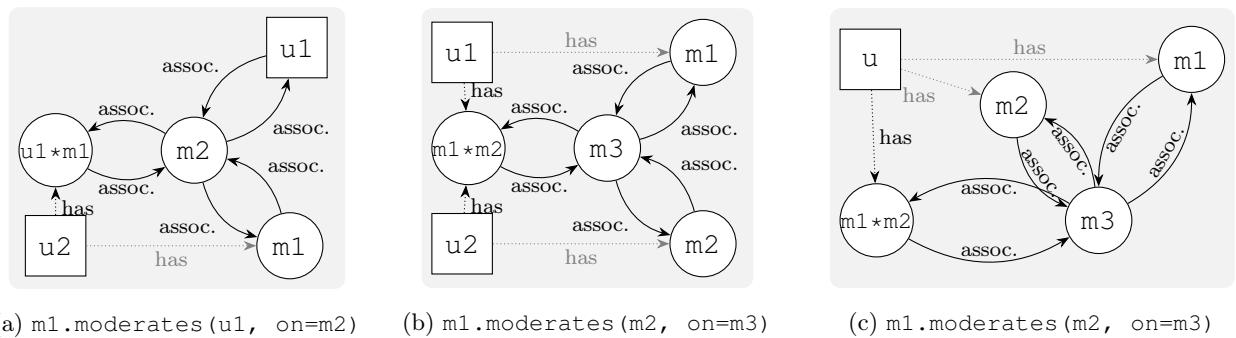


Figure B.1: More complex examples of `moderates` written in Tisane's study design specification language, and their representation in Tisane's graph IR. Variables are named with u for units, m for measures, and v for data variables that can be either units or measures. Black edges have been added due to the `moderates` relationship. Gray edges already existed in the graph. In (a), only $m1$ is a measure, whose unit is $u2$, so $u1*m1$ inherits an attribution edge only from $u2$. In (b), $m1$ and $m2$ are measures, with units $u1$ and $u2$ respectively, so $m1*m2$ inherits attribution edges from both $u1$ and $u2$. In (c), measures $m1$ and $m2$ share a unit, u , and $m1*m2$ inherits only one attribution edge from u .

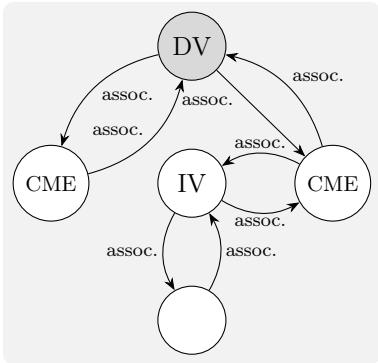


Figure B.2: A graph demonstrating an edge case for candidate main effect identification, where the graph contains only associative edges. Candidate main effects are labeled “CME”, independent variables “IV”, and dependent variables “DV”. Variables that are none of the above are left unlabeled. When a graph contains only associative edges, candidate main effects are identified as those that are either associated with the DV or are associated with both the IV and the DV. (Note that the graph could contain additional edges/nodes other than the ones pictured, but the additional edges would not violate any of the initial checks that Tisane makes on the graph IR.)

B.2 Cautioning analysts about adding certain kinds of variables

IVs
IVs: Interactions
Clustering
Data Distributions

Main Independent Variables

Main independent variables are variables whose influence on the dependent variable you are interested in.

Tisane derives main independent variables based on the `causes` and `associates_with` relationships you specify.

Warning: `age` suggested due to associative relationship

| | |
|--|---|
| <input checked="" type="checkbox"/> treatment <small>(i)</small> <input checked="" type="checkbox"/> motivation <small>(i)</small> <input type="checkbox"/> age <small>(i) Warning</small> | If <code>age</code> may actually be caused by the dependent variable, <code>pounds_lost</code> , adding this variable could invalidate your results. |
|--|---|

Continue

Proceed with caution!

Figure B.3: An example of the warning text given for potential confounding associations. When analysts hover over the “Warning” badge, a tooltip pops up that explains that they should be careful about adding this variable. Associative relationships may in actuality be causal relationships, and if in fact `pounds_lost caused age`, then adding `age` would invalidate the model.

Chapter C

Appendix C: Exploratory study materials

C.1 rTisane study materials

To be added when they are finalized