Tisane: Authoring Statistical Models via Formal Reasoning from Conceptual and Data Relationships

Eunice Jun

emjun@cs.washington.edu University of Washington Seattle, Washington, USA

Jeffrey Heer

jheer@cs.washington.edu University of Washington Seattle, Washington, USA

ABSTRACT

Proper statistical modeling incorporates domain theory about how concepts relate and details of how data were measured. However, data analysts currently lack tool support for recording and reasoning about domain assumptions, data collection, and modeling choices in an integrated manner, leading to mistakes that can compromise scientific validity. For instance, generalized linear mixedeffects models (GLMMs) help answer complex research questions, but omitting random effects impairs the generalizability of results. To address this need, we present Tisane, a mixed-initiative system for authoring generalized linear models with and without mixedeffects. Tisane introduces a study design specification language for expressing and asking questions about relationships between variables. Tisane contributes an *interactive compilation* process that represents relationships in a graph, infers candidate statistical models, and asks follow-up questions to disambiguate user queries to construct a valid model. In case studies with three researchers, we find that Tisane helps them focus on their goals and assumptions while avoiding past mistakes.

KEYWORDS

statistical analysis; linear modeling; end-user programming; enduser elicitation; domain-specific language; transparent statistics; validity

ACM Reference Format:

Eunice Jun, Audrey Seo, Jeffrey Heer, and René Just. 2022. Tisane: Authoring Statistical Models via Formal Reasoning from Conceptual and Data Relationships. In *CHI Conference on Human Factors in Computing Systems (CHI '22), April 29-May 5, 2022, New Orleans, LA, USA*. ACM, New York, NY, USA, 16 pages. https://doi.org/10.1145/3491102.3501888

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CHI '22, April 29-May 5, 2022, New Orleans, LA, USA
© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-9157-3/22/04...\$15.00

https://doi.org/10.1145/3491102.3501888

Audrey Seo alseo@cs.washington.edu University of Washington Seattle, Washington, USA

René Just rjust@cs.washington.edu University of Washington Seattle, Washington, USA

1 INTRODUCTION

Statistical models play a critical role in how people evaluate data and make decisions. Policy makers rely on models to track disease, inform health recommendations, and allocate resources. Scientists use models to develop, evaluate, and compare theories. Journalists report on new findings in science, which individuals use to make decisions that impact their nutrition, finances, and other aspects of their lives. Faulty statistical models can lead to spurious estimations of disease spread, findings that do not generalize or reproduce, and a misinformed public. The challenge in developing accurate statistical models lies not in a lack of access to mathematical tools, of which there are many (e.g., R [63], Python [52], SPSS [58], and SAS [24]), but in accurately applying them in conjunction with domain theory, data collection, and statistical knowledge [26, 38].

There is a mismatch between the interfaces existing statistical tools provide and the needs of analysts, especially those who have domain knowledge but lack deep statistical expertise (e.g., many researchers). Current tools separate reasoning about domain theory, study design, and statistical models, but analysts need to reason about all three together in order to author accurate models [26]. For example, consider a researcher developing statistical models of hospital expenditure to inform public policy. They collect data about individual hospitals within counties. Based on their domain knowledge, they know that counties have different demographics and that hospitals in these counties have different funding sources (private vs. public), all of which influence hospital spending. To model county-level and hospital-level attributes, the researcher may author a generalized linear mixed-effects model (GLMM) that accounts for clustering within counties. But which variables should they include? How do they account for this clustering? The three most common mistakes in modeling hierarchical data [9] lead to miscalibrated statistical power, "ecological fallacies" [49], and/or results that may not generalize, which impact not only the validity of research findings [3] but also enacted policies. How can the researcher avoid these issues?

To reduce threats to validity and improve analytical practices, how might we derive (initial) statistical models from knowledge about concepts and data collection? Inferring a statistical model raises two 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 present **Tisane**, a system for integrating conceptual relationships, data collection details, and modeling choices when specifying generalized linear models (GLMs) and generalized linear mixed-effects models (GLMMs). GLMs and GLMMs are meaningful targets because they are commonly used (e.g., in psychology [9, 35], social science [29], and medicine [3, 5]) yet are easy to misspecify for statistical experts and non-experts alike [3, 9]. 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.

Tisane provides a study design specification language for expressing relationships between variables. For example, the public health researcher can express that average county income is associated with hospital spending based on health economics theory or specify that hospitals exist within counties. Tisane compiles the explicitly stated relationships into an internal graph representation and then traverses the graph to infer candidate GLMs/GLMMs. In this process, Tisane engages analysts in interactive compilation. Analysts can query Tisane for a statistical model that explains a specific dependent variable from a set of independent variables. Based on the input query, Tisane asks analysts disambiguating questions to output a script for fitting a valid GLM/GLMM. Interactive compilation enables analysts to focus on their primary variables of interest as the system checks that analysts do not overlook relevant variables, such as potential confounders or data clustering that could compromise generalizability. Figure 1 provides an overview of this process.

To examine how Tisane affects real-world analyses, we conducted **case studies with three researchers**. The researchers described how Tisane focused them on their research goals, made them aware of domain assumptions, and helped them avoid past mistakes. Tisane even helped one researcher correct their model prior to submitting to the ACM Conference on Human Factors in Computing (CHI). These findings corroborate those from an earlier pilot study that informed our design process (see supplemental material).

We contribute (1) a study design specification language and graph representation for recording and reasoning about conceptual relationships between variables and data collection procedures (5), (2) an interactive compilation process that asks disambiguating questions and outputs code for fitting and visualizing a GLM/GLMM (6), and (3) three case studies with researchers that demonstrate the feasibility and benefit of prioritizing variable relationships to author linear models (7). We also provide an open-source Python implementation of Tisane. ¹

2 BACKGROUND AND RELATED WORK

We first provide brief background about data analysis practices, GLMs/GLMMs, and causal analysis. Then, we discuss how Tisane extends prior work on tools for conceptual reasoning, study design, and automated statistical analysis.

2.1 Data Analysis Practices

Studies with analysts have found that data analysis is an iterative process that involves data collection; cleaning and wrangling; and statistical testing and modeling [21, 31, 32]. To formalize their hypotheses as statistical model programs, analysts engage in a dual-search process involving refinements to their conceptual understanding and iterations on model implementations, under constraints of data and statistical knowledge [26]. Analysts incorporate and refine their domain knowledge, study design, statistical models, and computational instantiations of statistical models while creating statistical model programs. Tisane facilitates one formalization cycle in this iterative process: deriving statistical models from conceptual knowledge and data measurement specifications.

2.2 Generalized Linear Models and Generalized Linear Mixed-effects Models

Tisane supports two classes of models that are widely applicable to diverse domains and data collection settings [3, 5, 35]: Generalized Linear Models (GLMs) and Generalized Linear Mixed-effects Models (GLMMs). 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 $(5.2.2)^2$.

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).

2.3 Causal Analysis

There are multiple frameworks for reasoning about causality [44, 50]. One widespread approach is to use directed acyclic graphs (DAGs) to encode conditional dependencies between variables [20, 45, 56, 57]. If analysts can specify a formal causal graph, Pearl's "backdoor path criterion" [44, 46] 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 [60] due to uncertainty of empirical findings [61] and lack of guidance on which variables

 $^{^1\}mathrm{Tisane}$ is available for download on pip, a popular Python package manager. The source code is available at https://github.com/emjun/tisane.

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 [18]. When we do use these terms, we use the definitions from Kreft and De Leeuw [29].

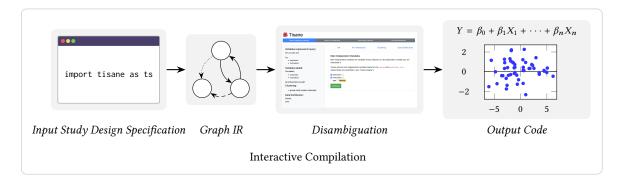


Figure 1: Overview of the Tisane system. Analysts specify a set of variable relationships (Input Study Design Specification). Tisane represents these in an internal graph (Graph IR). To infer a statistical model, Tisane engages analysts in an interactive compilation process that elicits additional input from analysts in a disambiguation process (Disambiguation) and outputs a script for fitting a valid GLM and visualizing its residuals (Output Code).

and relationships to include [67]. Accordingly, Tisane does not expect analysts to specify a formal causal graph. Instead, analysts can express causal relationships as well as "looser" association (not causal) relationships between variables in the study design specification language.

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 [56, 57]. Recently, VanderWeele proposed the "modified disjunctive cause criterion" [66] 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. Tisane 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 [66]. To guide analysts through the suggestions, Tisane provides analysts with explanations to aid their decision making during disambiguation.

Finally, GLMs 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.

2.4 Tools for Conceptual Reasoning and Study Design

Tools such as Daggity [64] support authoring, editing, and formally analyzing causal graphs through code and a visual editor. Daggity requires users to specify a formal causal graph, which may not always be possible [60, 61, 67]. Although a knowledgeable analyst could use Daggity to identify a set of variables that control for confounding to include in a linear model, Daggity does not provide this support directly. In contrast, Tisane aims to (i) help analysts may not be able to formally specify causal graphs and (ii) scaffold the derivation of GLMs and GLMMs from causal graphs.

Several domain-specific languages [1, 55] and tools specialize in experiment design [4, 15, 62]. A primary focus is to provide researchers low-level control over trial-level and randomization details. For example, JsPsych [13] gives researchers fine-grained control over the design and presentation of stimuli for online experiments. At a mid-level of abstraction, Touchstone [36] 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 [15] helps researchers design experiments based on statistical power. At a high-level of abstraction, edibble [62] 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. betweensubjects), 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.

2.5 Tools for Automated Statistical Analysis

Researchers have introduced tools that automate statistical analyses. Given a dataset, the Automatic Statistician [34] 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. In contrast, Tisane prioritizes analyst-specified conceptual and data measurement relationships and uses them to bootstrap the modeling process. As a result, Tisane aims to ensure that statistical analyses are not only technically correct but also conceptually correct.

AutoML tools automate machine learning for non-experts. Tools such as Auto-WEKA [65], auto-sklearn [16], and H2O AutoML [30] aim to make statistical methods more widely usable. Tisane differs

from AutoML efforts in its focus on analysts who prioritize explanation, not just prediction, such as researchers developing scientific theories. As a result, Tisane provides support for specifying GLMMs, which some prominent AutoML tools, such as auto-sklearn [16], omit. Tisane ensures that inferred statistical models respect expressed conceptual relationships. Thus, Tisane programs can serve a secondary purpose of recording and communicating conceptual and data measurement assumptions. In addition, Tisane explains its suggestions to users and guides them in answering disambiguation questions in hereas, 20110ML tools do not by default. Tisane's expla-

nations are grounded in the variable relationships analysts specify. Although H2O AutoML offers a model explainability module [14], which is a proper several analysts specify. Although H2O AutoML offers a model explainability module [14], which is a proper several analysts specify. These explainations take the form of plots without conceptual takes a graph anations take the form of plots without conceptual explainations. Take the form of plots without conceptual explainations take the form of plots without conceptual explainability. It is a several explaination of plots without conceptual explaination in the form of plots without conceptual explaination. The form of plots without conceptual explaination in the form of plots without an interest and explaination in the form of plots without an interest and plots are plots and plots and plots and plots are plots and plots and plots are plots and plots and plots and plots are plots and plots and plots are plots and plots and plots are plots and plots are plots and plots are plots and plots and plots are plots

a single statistical model, whereas Tea outputs a set of statistical tests.

Recent work in the database community helps researchers answer causal diestoches database community helps researchers and a Call for the province of the causal of the causal diestoches and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results. Like Carly Translets and a Call for show researchers results and a Call for show researchers. The show researchers are shown as well as the show researchers and the call of the call o

USAGE SCENARIO

3 USAGE SCENARIO
To illustrate how a researcher might use Tisane, we compare two By illineated Feeler Fine Philomich the sea Tie anna see, compose of any hamphatian lasa phala anno ina tha anno ata a toisis a en ira-(simplified from Labibly strates are the switters accorded tweeth Tit). general language lang Mish and raid and an hope at the area to provide a type effect states. newsart sine unninger then have elevel an edanity sinky latest Their resegrala ouse tiepais "il la myenus a dons des consiris e renei mese, affecte <u>weischt less? «The var reraittelvälle adultete ehe enst et Ansexernise</u> scoupenfiscuss fractics continuis is beneaus coment. The custom continuis garden harsignad dersanden germedregeimen en det besetber i 24 achunas an carngrium sutalu sapis ur eig. The ur engunder to promonent de l'hor edulteismentixationessessessessessessessessestatata beeinviinas ofitha. experiment and the interior and the residual and the interior and the contract of the contract Misharhwast on sandelys 147 herora teach teachainh Beidaat mea. Tiranei Mehilo bothiy, evernaricana deragarcher a femiliari veite theis. stracedifield of enudye Michael and Brislest our netatet oriestexactor Them have both used GLMs in the past but neither has heard of GLMMs.

Workflow in Python using statsmodels 3.1 Workflow in Python using statsmodels Michael takes a first attempt at creating a model. He loads the data Michealtakeen fiert attemproap casting omedel. Healneds the data

Michael uses statsmodels [47] to analyze the data³. Bridget uses Tisane. While both are experienced researchers, familiar with their shared field of study, Michael and Bridget are not statistics experts. They have both used GLMs in the past but neither has heard of GLMMs.

Workflow in Python using statsmodels 3.1

Michael takes a first attempt at creating a model. He loads the ables. The first model (Listing 1) that Michael tries has the dependent ivertable the independent at the lindependent controller. pregional quantities (independent extriable) are given visconorol vs.

```
treatment) and motivation.
import statsmodels.formula.api as sm
import statsmodels.formula.api as sm
23 import pandas as pd
45 data = pd.read_csv("data.csv")
 data['regimen_condition'] = data['regimen_condition'].
      astype('str').astype("category")
 data['group'] = data['group'].astype('str').astype("
       category")
 m1 = sm.glm("pounds_lost ~ regimen_condition + motivation
  ", data=data)
Listing 1. whenacts mist model attempt, from the usage
```

Listingid: (Michael's). funt medebattemptourem 1 be usage steparite (Section 63); Mighael specifical pautosal pot as his danendaria bleand regimen and motivation as his

independent variables.
Although this model includes the primary variable of offers a rindunister i kontrologora dike belahari tiha ga sang sagad terat Health and the three three will be also also the contract that the action of the contract that the action of the contract three thre encined in a rational participation of the property of the pro ishidity is not a the and directly species on a species that can difficult the amaquiether litely legit in the design that are difficult to measure and phis school model (existing 2), Michael adds group as an additismah independent variable: Generaleting ahat dhis model accounts tion variables dentaining to individual adults his moderns clost, motivation and incompanitivities in the second in the second seco and series ship where supplies the state of the series and series are series and series and series and series are series are series model is good enough and accepts it as his final model.

 $\label{eq:m2} \begin{array}{ll} m2 = sm.glm("pounds_lost ~ regimen_condition + motivation \\ m2 = sm.glm("pounds_lost ~ regimen_condition + motivation \\ \end{array}$

+ group", data=data)
Listing 2: whenael's second model attempt. Dunding off Listing se: Michaelis gregord, madel attemet a Building on a historist den del richieting out. Michael adds an additional independent variable, group.

Workflow with Tisane

3.2. Workflow with Tisane Bridget starts by listing the variables of interest. Lines 3-8 in List-Brijdest stattskividitinie elekariabled of inferratationes inesinasion itris fush own having the rial telepolar deviated and the same at home as the thereolymphine electron that he was pointed to the conservation of dankingsintsy i Honoverna does en telefolio en en en telefonia franchi i kent i france de en en en en entelesc most identification to the control of the control o venesei levinius i explitere eletti aste pennyle This Tiral une vuseis er en alces the specific line xplicit are cause Bridget has data, she does not need

```
to specify the cardinality of variables. import tisane as ts
 The workflowaineR is almost identical to the workflow in Python using statsmodels.
34 # Variable declarations
45 adult = ts.Unit("member")
56 motivation = adult.numeric("motivation")
g pounds_lost = adult.numeric("pounds_lost")
78 group = ts.Unit("group")
 regimen = group.nominal("regimen_condition") # control vs
```

Listing 3. The mot simplet of the example Hoane program,

ha firetieminnat of

statistical model in line 15.

3.2. Workflow with Tisane
Query Tisane for a statistical model

the data Γhe first variable

tsmodels.

```
er in R).
ects of a
ss. Their
en affect
exercise
earchers
other 24
ured the
ig of the
eriment.
get uses
ith their
experts.
heard of
```

a disambiguation process (see Figure 2) to generate a final output statistical modeling script. adult = ts.Unit("member") motisane launghes a filmlinthembrowser after executing the progrant. In the Golf Tisano asks (Bridgets to look) over her choice of \$apiæblest(Figuite(2)βrAspseen in panel B, there are no additional variables to consider. Bridget continues to the next tab, interaction effects (panel D). Tisane explains that interaction effects do not Listing 3: The first snippet of the exemple this presping ramspecifies by her by her ved y aria bles ane has automatically included

design = ts.Design(dv=pounds_lost, ivs=[regimen_condition -

ts.infer_model(design=design)

Visting 5nThe final snippet of the example Tisane program.

Bridgethquariesi Tisane doe grotatis tinal madelphyspecifyikes

Thee Ttrebel designers, modern imparation has rightly execute ordering a

blas plad Pashed dons it had in our played piece para var dinnisis to figa of a billiois

listing as well as Listing 3 and Listing 4, Bridget engages in

usane: Authoring Statistical Models via Formal Reasoning from Conceptual and Data Rela

, motivation]).assign_data("data.csv")

By expressing how the variables relate to one another. Bridget ex**By expressing how the variables relate to one another Bridget** ex by expressing now me variables plate to one constitute, bridge becomes more consciously aware of the assumptions the research becomes more consciously aware of the assumptions the research team has made about their domain. In the 10.1 sting 4. Bridge team has made about their domain in the 10.5 sting the research ream has made about their domain in the 10.5 sting the first process pecules that regimen directly causes pounds lost while in the specules that regimen directly causes pounds lost while in the specules that regimen directly causes pounds lost while in the specules that regimen directly causes pounds lost while in the specules that regimen directly causes pounds. The process of the process in street Bridget an Swers questions about the the dependent variable pounds less to identify family and link functions. She specifies That the dependent variable is continuous. She chooses the Gaussian adult nests within (group) tamily with the detault link function. Finally, Bridget clicks on the Listing offen Argentinustione of whoe stripped time is ting in the control of the stripped time in the control of the stripped time in the control of the co

*As explained in subsection 6.3, if Bridget had executed the script in a Jupyter notebook, the New will have the himselfer. By intercept expressions a study except the New will have the himselfer. By intercept expression of a study expression of the study of the himselfer. By the himselfer in t

```
# Query Tisane for a statistical model
```

ts.infer_model(design=design)
Listing 5: The final snippet of the example Tisane program, Listing 5: The final snippet of the example Tisane program, Listing 5: The final snippet of the example Tisane program, continuing from Listing 1: Listing 1: The final snippet of the example Tisane program, continuing from Listing 1: Listing 1: Listing 1: Listing 1: Listing 1: Listing 2: Listing 1: Listing 2: Listing 3: Listing 3: Listing 3: Listing 4: Listing 3: Listing 4: Listin

output statistical modeling script.

Tisane launches a Giff in the browser after executing the program of the p

lysts 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 (Listing 4), Tisane infers candidate main effects that may be potential confounders. Tisane asks analysts if they imond like to include the se variable sy explaining imas tooltip (C) why the variable may be important to include. (D) Tisane only suggests interaction effects if analysts specify moder-(E.) Bridget sees that Tisane has automatically included exercise ating relationships in their specification. This way, Tisane group as a random intercept. (Isane does not include a random ensures that model structures are conceptually justifiable. slope for group because there is only one observation fer adult in (E) From the data measurement relationships analysts programment see to 2.3. Bridget has not heard of a random intercept beyonde (line 15 in Listing 4). Tisane automatically infers and includes readom effects to increase generalizability and exfor exercise groups is necessary since adults were in groups and ternal validity of statistical findings. (E) Tisane assists anaternal validity of statistical findings (F) Tisane assists and answers questions about the the dependent variable pounds 10st to identify family and link functions. She specifies that the dependent variable is continuous or about count data?). To help ana-dent variable is continuous. She chooses the Gaussian family with the default link function. Finally, Bridget clicks on the button to generate code. Fisane's output script contains code to fit the stawritten in the study design specification language from the generate code. Pisane's output script contains code to fit the statistical works and works with the statistical works are authority statistical works are all the statistical works are al

Bridget author the following GLMM:

```
model = Lmer(
   formula="pounds_lost ~ regimen_condition +
motivation + (1|group)",
   family="gaussian",
```

Listing 6 Bridget's output statistical model from using Tisane. Tisane suggests a GLMM with group as a random intercept. The output script contains code for fitting his model and inspecting it using a residuals plot.

3.3 Key differences in workflows and statistical f you believe you omitted a modificing relationship, go back to your program and specify it saring the minimal to the control of the control

Even with experience mode og in Python, Michael makes two common mistakes in authoring inear models [9, 29]: disaggregating

Bidurelee V. shiffiereficesei (II. W to kallows and ista fistical-figure 2: Example I isane GUI for disambiguation from usage scenariotal isane asks analysts disambiguating questions about wariables that are conceptually relevant and that analysts disambiguating duestions about wariables that are conceptually relevant and that analysts may have overlooked in their query. (A) the left hand 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 analyst structings. (B) based on the variable relationships analyst specify (I isling 4). I sake inters candidate main effects that may be potential confounders. It same asks analyst sit they would like to include these variables explaining in a footing to what the variable may be important to include these variables explaining in a footing to what the variable may be important to include these variables explaining in a footing to what the variable may be important to medical distance only suggests interaction effects that model structures are conceptually listing and the same only suggests interaction effects that onclude the same and the same asks analysts if they want the variable may be important to include (B) disante only suggests interaction effects that onclude (B) disante only suggests interaction effects that onclude the same and the and external validity, avoiding errors and reducing the cognitive burden along the way. Once fit, Michael's and Bridget's final models (see) disagreeon the precise effect sizes of regiment condition and motivation, which are partinent to their motivating research question. Michael concludes that regimen_condition and motivation are less important than they really are. Additionally, the coefficients and standard errors for each group suggest that Michael's model overlooks important group differences. Therefore, using a GLM instead of formula="pounds_lost ~ regimen_condition +

```
motivation + (1|group)",
   family="gaussian",
   data=df,
```

Listing 6: Bridget's output statistical model from using fisance. Tisance suggests a GLMM with group as a random tisance. Tisance suggests a GLMM with group as a random intercept. The output script contains code for fitting this metacent are output script contains code for fitting this

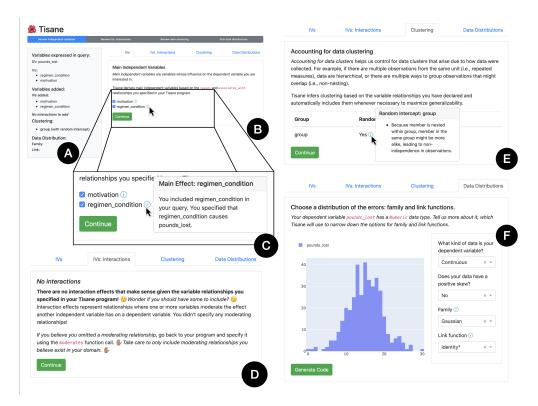


Figure 2: Example Tisane GUI for disambiguation from usage scenario. Tisane asks analysts disambiguating questions about variables that are conceptually relevant and that analysts may have overlooked in their query. (A) The left hand panel gives an overview of the model the analyst is constructing. (B) Based on the variable relationships analysts specify (Listing 4), 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 (line 15 in Listing 4), 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.

the appropriate model, a GLMM, leads Michael to answer the research question differently than Bridget, artificially inflates statistical power [9], and compromises the generalizability of his findings $[3].^5$

4 DESIGN GOALS

We articulated four design goals based on prior research and our formative work. The supplemental material details our design process.

DG1 - Prioritize conceptual knowledge. Current tools require analysts to transition back and forth between their conceptual concerns and their statistical model specifications using math and/or code [26]. Analysts' conceptual knowledge remains implicit and

	$eta_{ ext{regimen_condition}}$	p	$eta_{ exttt{motivation}}$	p	$eta_{ t group}$	p
Michael	1.628	.046	3.119	.000	Var.	Var.
Bridget	1.659	.005	3.193	.000	N/A	N/A

Table 1: The coefficients for each of the independent variables in Michael's and Bridget's models. "Var." stands for "Various," since there were multiple coefficients generated. The complete output tables for Michael's and Bridget's models are included in supplemental material.

hidden [38]. As a result, analysts may resort to familiar but suboptimal statistical methods [26] or accidentally overlook details that lead to conceptually inaccurate statistical models. One solution is to provide tools at a higher level of abstraction that allow analysts to express their conceptual knowledge directly. However, a higher level of abstraction alone is not enough. Tools must then leverage the expressed conceptual knowledge to guide analysis authoring.

⁵Another common mistake, not shown here, is to aggregate observations and use group means of the independent variables in the model, artificially deflating statistical power ("ecological fallacy" [49]). Kreft and De Leeuw [29] share an example where disaggregating vs. aggregating data lead to different signs for a fitted parameter. Unfortunately, we could not access the data to illustrate this here.

Tisane provides a high-level study design specification language that captures the motivation behind a study (Section 5). Tisane represents the specification in an internal graph representation to derive only conceptually accurate statistical model candidates. To arrive at an output statistical model, Tisane asks analysts disambiguating questions and provides them with suggestions and explanations based on their expressed variable relationships (Section 6). Importantly, Tisane does not fit or show modeling results during disambiguation to discourage statistical fishing. Although Tisane does not prevent researchers from re-starting and iterating on their Tisane program to attain specific statistical model findings, Tisane programs act as documentation for conceptual relationships that others could audit.

DG2 - Prioritize the validity of models. At present, the burden of valid statistics lies entirely on analysts. Tisane divides some of this burden by (i) ensuring correct application of methods (i.e., GLM vs. GLMM) and (ii) inferring models that increase the generalizability of results for GLMMs [2, 3]. Tisane helps analysts author GLMs and GLMMs that satisfy two assumptions: (i) observational dependencies and (ii) correct family and link functions. First, Tisane infers and constructs maximal random effects that account for dependencies due to repeated measures, hierarchical data, and non-nesting compositions. Maximal effects structures account for within-sample variability and thereby mitigate threats to external validity due to sampling biases from the choice of observational units and settings [54]. Second, Tisane narrows the set of viable family and link functions to match the dependent variable's data type (e.g., numeric). Tisane's GUI asks follow-up questions to determine the semantic type of variables (e.g., counts), further narrowing analysts' family and link function choices. The output script also plots model residuals against fitted values and provides tips (as comments) for interpreting the plot. The family and link functions Tisane suggests are intended to bootstrap an initial statistical model that analysts can examine and, if necessary, revise. This is how Tisane helps analysts avoid four common threats to statistical conclusion and external validity [11]: (i) violation of statistical method assumptions, (ii) fishing for statistical results, (iii) not accounting for the influence of specific units, and (iv) overlooking the influence of data collection procedures on outcomes.

DG3 - Give analysts guidance and control. Analysts may have insight into their research questions and domain that a system cannot capture. At the same time, analysts, especially those with less statistical experience, may lack the knowledge to select among many possible statistical models, which may inadvertently encourage cherry-picking based on observed results. Thus, Tisane adopts an interaction model that asks analysts specific questions to resolve modeling ambiguity rather than show multiple statistical models at the same time. Tisane also does not automatically select a a "best" model (e.g., highest R², easiest to interpret) but rather gives analysts suggestions and explanations to help them come to a statistical model that is valid and appropriate for their goals.

DG4-Facilitate statistical planning without data. Experimental design best practices, such as pre-registration, encourage researchers to plan their statistical analyses prior to data collection. Tisane supports these best practices by not requiring that analysts provide data. If analysts do not have data, analysts must specify the cardinality at variable declaration. Without data, Tisane cannot

validate variable declarations, but in this case, Tisane still guides analysts through the same interactive compilation process. The output Tisane script will include an empty file path and a comment directing analysts to specify the path to their data prior to execution. Analysts could attach this output Tisane script to their pre-registrations. After analysts collect data, they can re-run their previously specified Tisane program to validate and inspect their data. If Tisane does not issue any validation errors, analysts can proceed to execute their script.

5 STUDY DESIGN SPECIFICATION LANGUAGE AND GRAPH REPRESENTATION

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.

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.

5.1 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.

5.2 Relationships between Variables

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 3). 4 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.

5.2.1 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 [60, 61, 67], prior work has observed that researchers already use informal graphs that contain associative relationships when reasoning about their hypotheses and analyses [26], 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 [23], 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 (3(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 (3(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 (3(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.⁶

- 5.2.2 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 [8, 19, 29]:
 - **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) [19].

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.

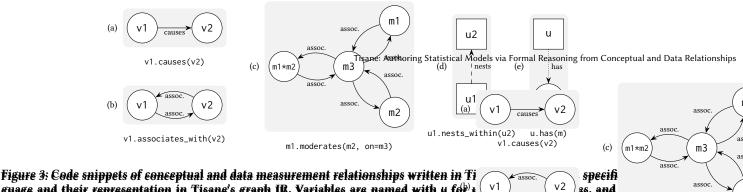
Thus, Tisane addresses the three problems by (i) eliciting data measurement relationships from analysts to infer clusters and (ii) formulating the maximal random effects structure, optimizing for external validity (6). Below, we describe language features for expressing data measurement relationships.

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 (3(d)). There is one edge from the nested unit (e.g., adult) to the nesting unit (e.g., group) ⁷.

Frequency of measures: Repeated measures, Non-nesting composition. When a measure is declared through a unit, Tisane adds an

⁶In statistical terms, this means that within-level interactions have one unit while cross-level interactions may have two or more units.

⁷The GitHub repo contains a gallery of examples that include nesting relationships.



guage and their representation in Tisane's graph IR. Variables are named with u for variables that can be either units or measures. All edges depicted are those that are according to the relationsl moderates example, we assume that m1 and m2 both belong to the same unit, and for simplicity, the attribution edge (labeled as "has") from m1 and m2's unit is not shown. For more complex examples, see the supplemental materials.

attribution edge ("has') from a unit node to a measure node (3(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 [19] composition arises when necessities 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 observes 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 obregyation of the and prenepriolate the assumption of interned energing that ships frake the usage scenario. causes edges are labeled with headspalysts declared with earlies are lifeted from each soch Dassieure agethinusch the neutre with instances napameter doison cargeter accente nasintenacionariable, a Tisane Exactíy 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 in ough the per function. For INTERACTUALLY QUERVING THE GRAPH on assig**iR**d: device = subject.nominal('Input AMBPespecifying tangaste teletratish bis? Indevise caditien) it is the nerafunctionalses the Timerevarie bleast ardinality special all but cap inestgad wan a dataeyariah le'ar nunberhe Galine ta be es eku cangifying secorandipadity-talses parameter to aner pMeesee the of city ingiagneesures's a vinterical-inatences to phancinterpie confaction

eugee for a sing the fixactal xi one into its one cibiling, a controlle is ever

tacting arefor expressing and Family and link functions (9) input

elicToadstetuvinsathsigteseeps ssturpeatedin, east (4) generations ting

composition: Tisana computes the number and the stances of more

surps. and their are lations binctore there exercises. Whas are surpat

preglachted with a new data of this target percentage are non-

eidened, the vary this pieces unit to Measures of the declared with

Figure 3: Code snippets of conceptual and data measurement relationshing the conceptual and data measurement relationshing the content of 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 technology and the conceptual have little to no conceptual for the checks pass, Tisane pounds_lost unsyt phase.

6.2 Candidate stati odel 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.

Figure in The graph representation of the variables and refationships from the charge scenario: causes edges are labeled with causes. associates with edges are labeled with charge ice. District edges indicate nests within relationships, and dotted edges indicate has remaining pain 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 the say STATISTICAL MODEL INFERENCE:

62.15 TATES COMENS PER PROPERTY OF THE STATE OF THE ACTION OF THE PROPERTY AND THE PROPERTY OF THE PROPERTY AND THE PROPERTY OF THE PROPERTY O

6.1]

At the binput st correctr study de variable any of ticause froissues a Tisane chave litt. If the ch

A GLM/function

of main, concept effects a fects. To type of Tisane g

disambi The pones an erroneo pre-emp may aris

6.2.1 L tistical r set of or the que addition effects in statistical disjunct

⁸Tisane cu identify co

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*.

6.1 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.

6.2 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.

- 6.2.1 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 6.1, 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 model⁸. The first two sets below are from the "modified disjunctive cause criterion" [66]:
 - Causal parents. For each IV in the query, Tisane finds its causal parents (see 5(a)).
 - Possible causal omissions. Tisane looks to see if any other variables not included as IVs cause the DV (see in 5(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(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 (see Figure 3 in supplemental material).

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.

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.

- 6.2.2 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
- 6.2.3 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 Section 5.2.2). 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 [2, 3].

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 [2]). Tisane handles correlation of random slopes and intercepts during disambiguation (section 6.3). Maximal random effects may lead to model convergence issues that

⁸Tisane currently treats each input IV as a separate "exposure" variable for which to identify confounders. Tisane then combines all confounders into one statistical model.

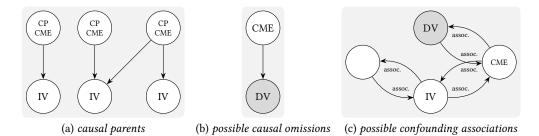


Figure 5: Graphs demonstrating causal parents, possible causal omissions, and possible confounding associations. In graphs (a) and (b) (left and middle), all edges are causal. Independent variables are marked "IV", discovered candidate main effects "CME", dependent variables "DV", and causal parents "CP".

analyategoddrany by letactromovinegus; pdeling indennadeut variablesand mandern reffective Navior thanks inter time on it bounds implumed in the control of the cont (nodelic important the voted grip of that future privisions are also valid (DG2 thealidity) ustering tab, Tisane shows analysts which random effects it automatically includes based on the selected main and in the factorial of the type determines the set of candidates amily and other types the first the set of candidates amily and the first types the first type the set of example and the control of the contr distributions c Similarly inominally asiables are 19st all owed to have Ganssian distributions a Few thermore each foreign date that of nonsible link functions. For example, a Gaussian family distribution mexshaye or Identity a hoganu Sounton) Poot dink dunctione Therstatiss tics literature idecuments possible combinations of family and link functions for appeaific data day the ps analysts examine their data and isenee includes are ramny family distribution to a cyandidate famisereand their applicable link functions be its current implamentation. Titapending on state snooteles (1531-1601-161-161 and degrees 141/25). foreGIAM on Assuch Tisane is limited to the family and link function pairings implemented in these libraries. As statsmodels' and pyrer 4's ishinpertatistical Modern wonsither, friture naisones can be **EXTERNALLY** functions considered based on the declared data type of variables (see 6.2.4) and lightweight viability checks, such as ensur-6.3 tha Eliciting denial wish Input afort Disambigutation The disambigative integers alvested an opportunity for a day start explore the space or generated models was edular their of ganal quety. Given our design considerations to prioritize conceptual knowledge (DG11-s Contestual knowledge) have give analysts guidance (DG3 = Cultaince anarcontrol, we designed a GUT to scaffold analysis reasoning and elicit their thout. For versatility, we implemented Tisane's GOT using Plotly Dash [10]. Analysis can enther execute their Tisane programs and use the GUI inside a Jupyter notebook (no additional widgets needed) or run their Tisane programs in an DE or terminal, in which case Tisane will open the GUI in a web Browser Figure 2 tgives tan inversitive of the GUI on: (i) a log of GUI ch Candidate istatistical models are organized according to (i) independent variables (main effects and interaction effects); (ii) data chistering (rathdom effects); and (tii) idata distribution (family and link functions) In the main effects tab; Tisane asks analysts if they worlde diktorer saciatifizationalides abstitute inthinvendel sandaploit sasielialoragai bet datte eptalales in levdent. Una de sante capiprophieteneds, Tisane suggests moderating relationships to include but does not 'See supplemental material for a complete listing of Tisane's supported family and automatically ring lude them because analysts may not have specific

by pathity enit wilk in interactions (Direct Wideness and ing walkily analysis denestors: The any parderaping solutions bigs a climated explaining whan to literate in the characteristic and the first of the consequence of

The ignilytaly. Intradistribution, thelps and year by annier lost relate and selicifical in that conducting land and yeight to try. Appropriate selection of the appendies recognitive international when we have allowed the dependence of the appendies recognition and the alien by the above the independence of the appendies recognition and the alien when the alien and the alien alien and the alien and the alien and the alien and the alien alien and the alien and the alien alien and the alien alien and the alien alien and the alien alien alien and the alien alien alien and the alien alien

For assimatial sangsing futured at rojects identified what odd it is not of family post ions escardered was from this declared data type of variobles (sen 62-4) and light weight wishlitty the cks agusto as census inaithatal Poissons diverbution on any vannicable sfor variable athat bayetagen og attive integer proless. Tisane asksing eations designed town cover, more segment call ver can ingly lateraty presidence. Coverts) there are incovinced at variable declaration evanaly strivith out data can mass were 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 stat smodels [47] and pymer 4 [25]. Analysts can then choose a different link function, as long as it is supported to work indendently. We answered any questions researchers had while using 614 ane Each study session lasted approximately 2 hours. At the The revocof who other we seather in the Raction of the plan (iii) dato exce Tistable again debing the pleant dwii) madnihof GUI choices. To increase transparency of the authoring process, Tisane provides a log of 75.4 selection Study I! Planning at one which they analyst can include in pre-registrations for example (DG4 - Statistical planning). RI, a clinical psychology PhD student, had recently submitted a paper and was planning a follower psRdgreported statished had never takerrationnal class on modeling techniques but taught herself for

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.

7 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.

7.1 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+ [39] and SPSS [58]. 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 [43] for R1 that previous tools had not.

7.2 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 [60, 61, 67] 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.

7.3 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).

7.4 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 subsection 6.3.

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).

8 DISCUSSION

In this work, our motivating question was "How might we derive (initial) statistical models from knowledge about concepts and data collected?" This question presented two challenges: (i) how to elicit the information necessary to author a GLM/GLMM and (ii) how to computationally infer a valid statistical model given this information. To address the first challenge, we designed and developed Tisane's study design specification language. To address the second challenge, we developed a graph representation that Tisane traverses to derive candidate statistical models. We also developed a novel interaction model that involves interactive compilation to address both challenges. Throughout the design process, we employed statistical methods and theory, theories of how people analyze data, and an iterative design process with researchers. When using Tisane, researchers in our case studies reported focusing more on their analysis goals and becoming more aware of their assumptions and even identified and avoided previous analysis mistakes. Below, we reflect on future opportunities for Tisane to further enhance statistical practice, interpretation of results, and the end-to-end data analysis pipeline.

Design for statistical validity. Campbell's theory of validity – encompassing statistical conclusion, internal, external, and construct validity [6, 11] – has influenced disciplines widely (e.g., [54]), including epidemiology (e.g., [37]), software engineering (e.g., [42]), and psychology (e.g., [6]). Viewed through the Campbellian framework, Tisane helps analysts avoid four common threats to statistical conclusion and external validity: (i) violation of statistical method assumptions, (ii) fishing for statistical results, (iii) not accounting for the influence of specific units, and (iv) overlooking the influence of data collection procedures on outcomes [11].

Tisane fills a need to align analysts' conceptual models with the statistical models they want to implement but find difficult to express with the current tools available. By integrating conceptual, data, and statistical concerns, Tisane facilitates the hypothesis formalization [26] process, which can be an error-prone and cognitively demanding process that existing tools do not yet support.

In the future, we plan to develop additional strategies for enhancing the validity of analyses authored with Tisane. As discussed in Section 6.3, our current approach to family and link functions is only an initial step. We look forward to developing and comparing multiple strategies for scaffolding the family and link function selection and revision process. For example, what if the Tisane GUI allowed analysts to fit multiple models that varied in their family and link functions, plotted each model's residuals against the predicted values, and gave analysts visual guides for comparing models? To avoid false discovery rate inflation, Tisane could partition analysts' data, fit models to only a subset, and output a script for fitting a selected model using another subset. Although possible for large datasets, this strategy would encounter limited statistical

power for smaller datasets. Alternatively, what if Tisane calculated Bayes factors for variables in the models [12, 17, 48] after analysts tried multiple family and link combinations? Carefully balancing statistical rigor and usefulness to domain researchers who may be statistical non-experts deserves careful consideration.

Prevent p-hacking. Tisane generates a space of possible models from a set of conceptual and data measurement relationships. By querying Tisane for a model, analysts will only consider a set of models that are compatible with these relationships. As a result, Tisane helps analysts avoid unintentional p-hacking. Especially motivated p-hackers could specify questionable conceptual and data measurement relationships to manipulate the space of models Tisane generates. However, in this case, review or inspection of the Tisane program during pre-registration or peer-review, for example, could identify such malicious practices. In these ways, we believe that p-hacking is more difficult in Tisane than in existing analysis tools.

Future work to further discourage p-hacking could extend Tisane to conduct a sensitivity analysis on the space of possible models and only report models and results that are robust across the space. A challenge in this approach is that statistical non-experts may need more scaffolding to understand and interpret the results of sensitivity analyses.

Scaffold interpretation of statistical results. Tisane's focus is on authoring GLMs/GLMMs, but accurate interpretation is also necessary. For instance, analysts may need help interpreting what their statistical models and results mean in relation to their input conceptual models. Do the results suggest their conceptual model is correct? What kind of inferences should they make? Future work should address these interpretation challenges, which may require eliciting hypotheses and expected results from analysts.

Although researchers in our pilot or case studies did not presume Tisane helped with formal causal analysis, the ability to express causal relationships (causes) may lead some analysts to erroneously assume that their models assess causality. Changing the name of the language construct and/or building out support to interpret GLM/GLMM results may resolve this concern. One way to support accurate interpretation and reporting could be to output a figure representing the input conceptual model along with visual summaries of the data and/or statistical model for direct inclusion in publications. Tisane could also allow analysts to annotate their disambiguation decisions with their own rationale and provide a richer log of selections than currently supported. Tisane could even accept these augmented logs to save the state of the GUI in between analysis sessions.

Provide discipline-specific language support. When designing Tisane's study design specification language, we analyzed and developed language constructs common across existing libraries for study design (see supplemental material). In our case studies, we found that researchers had different conventions for describing their data ("unbalanced" (R2) vs. always "square" (R3)) and models (e.g., "controls" (R1) vs. "covariates" (R3)). This observation suggests opportunities to increase the usability of Tisane's SDSL by

providing syntactic sugar that may be more familiar to users. In the future, we plan to formally assess usability and identify "natural" programming [40] constructs that differ across disciplines.

An additional strategy for supporting more discipline-specific programming models and analysis needs is to integrate Tisane with existing study design libraries. For example, HCI researchers may find the lower-level randomization details that Touchstone2's interface [15] provides more natural. A system could summarize these details into the higher-level data measurement relationships in Tisane to bootstrap interactive compilation and output a possible statistical model. In this way, Tisane's graph IR can provide a "shared representation" [22] between study design tools and Tisane.

Integrate into end-to-end analysis workflows. Researchers in our case study were more comfortable with R. R1 and R3 expressed it could be helpful to have Tisane in R as an RStudio plug-in, for example, to fit into their workflows. As more users adopt Tisane, we will add an implementation in R.

Moreover, analysts may need to add or remove variables from Tisane's generated statistical models in order to accommodate model convergence failures, new data, or changing domain knowledge. However, adding or removing variables may subtly change the hypotheses analysts can test statistically. We look forward to extending Tisane to support model iteration, which presents two technical challenges: (i) recognizing when conceptual revisions are necessary and (ii) identifying and suggesting model changes that maintain conceptual validity or, at the very least, quantify conceptual shifts. Furthermore, in the model revision process, analysts may consider multiple alternatives. As R3 described, he preferred to run multiple variations of a model and compare them, a workflow akin to a multiverse analysis [59]. Given that Tisane already generates a combinatorial space of candidate statistical models, Tisane could generate a multiverse script for Boba [33] instead. A multiverse could help check the robustness of findings, and Boba's visual analyzer could help analysts further develop an understanding of their data and modeling choices. A multiverse may also help analysts explore and compare family and link combinations as well.

Tisane is one tool designed to enable analysts with limited statistical expertise to author valid statistical models. Tisane enables future possibilities and raises open research questions for creating an ecosystem of analysis tools that align tool interfaces with analysts' conceptual goals.

ACKNOWLEDGMENTS

We thank Yang Liu and Younghoon Kim for early feedback on Tisane's API; Leilani Battle, Matthew Conlen, Sherry Wu, and Rock Pang for feedback on early drafts of this paper; Tyler McCormick for feedback on the project and paper; and Maureen Daum for insightful conversations about how Tisane's graph IR relates to data models. We also thank the anonymous reviewers who provided valuable feedback.

REFERENCES

- Eytan Bakshy, Dean Eckles, and Michael S Bernstein. 2014. Designing and deploying online field experiments. In Proceedings of the 23rd international conference on World wide web. ACM, 283–292.
- [2] Dale J Barr. 2013. Random effects structure for testing interactions in linear mixed-effects models. Frontiers in psychology 4 (2013), 328.

- [3] Dale J Barr, Roger Levy, Christoph Scheepers, and Harry J Tily. 2013. Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal* of memory and language 68, 3 (2013), 255–278.
- [4] Graeme Blair, Jasper Cooper, Alexander Coppock, and Macartan Humphreys. 2019. Declaring and diagnosing research designs. American Political Science Review 113, 3 (2019), 838–859.
- [5] Benjamin M Bolker, Mollie E Brooks, Connie J Clark, Shane W Geange, John R Poulsen, M Henry H Stevens, and Jada-Simone S White. 2009. Generalized linear mixed models: a practical guide for ecology and evolution. Trends in ecology & evolution 24, 3 (2009), 127–135.
- [6] Donald T Campbell and Julian C Stanley. 2015. Experimental and quasiexperimental designs for research. Ravenio Books.
- [7] Herbert H Clark. 1973. The language-as-fixed-effect fallacy: A critique of language statistics in psychological research. Journal of verbal learning and verbal behavior 12, 4 (1973), 335–359.
- [8] Jacob Cohen. 1988. Statistical power analysis for the social sciences. (1988).
- [9] Jacob Cohen, Patricia Cohen, Stephen G West, and Leona S Aiken. 2013. Applied multiple regression/correlation analysis for the behavioral sciences. Routledge.
- [10] Plotly Dash Community. [n. d.]. Plotly Dash. https://plotly.com/dash/
- [11] Thomas D Cook, Donald Thomas Campbell, and William Shadish. 2002. Experimental and quasi-experimental designs for generalized causal inference. Houghton Mifflin Boston, MA.
- [12] Claudia Czado and Adrian E Raftery. 2006. Choosing the link function and accounting for link uncertainty in generalized linear models using Bayes factors. Statistical Papers 47, 3 (2006), 419–442.
- [13] Joshua R De Leeuw. 2015. jsPsych: A JavaScript library for creating behavioral experiments in a Web browser. Behavior research methods 47, 1 (2015), 1–12.
- [14] H2OAutoML developers. 2021. H2OAutoML v3.32.1.6 Documentation: Model Explainability. (2021). https://docs.h2o.ai/h2o/latest-stable/h2o-docs/explain. html#model-explainability
- [15] Alexander Eiselmayer, Chatchavan Wacharamanotham, Michel Beaudouin-Lafon, and Wendy Mackay. 2019. Touchstone2: An Interactive Environment for Exploring Trade-offs in HCI Experiment Design. (2019).
- [16] Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. 2015. Efficient and Robust Automated Machine Learning. In Advances in Neural Information Processing Systems, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (Eds.), Vol. 28. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2015/file/ 11d0e6287202fced83f79975ec59a3a6-Paper.pdf
- [17] Alan E Gelfand and Dipak K Dey. 1994. Bayesian model choice: asymptotics and exact calculations. *Journal of the Royal Statistical Society: Series B (Methodological)* 56, 3 (1994), 501–514.
- [18] Andrew Gelman. 2005. Why I don't use the term "fixed and random effects". https://statmodeling.stat.columbia.edu/2005/01/25/why_i_dont_use/
- [19] Andrew Gelman and Jennifer Hill. 2006. Data analysis using regression and multilevel/hierarchical models. Cambridge university press.
- [20] Sander Greenland, Judea Pearl, and James M Robins. 1999. Causal diagrams for epidemiologic research. *Epidemiology* (1999), 37–48.
- [21] Garrett Grolemund and Hadley Wickham. 2014. A cognitive interpretation of data analysis. *International Statistical Review* 82, 2 (2014), 184–204.
- [22] Jeffrey Heer. 2019. Agency plus automation: Designing artificial intelligence into interactive systems. Proceedings of the National Academy of Sciences 116, 6 (2019), 1844–1850.
- [23] Rick H Hoyle. 1995. Structural equation modeling: Concepts, issues, and applications. Sage.
- [24] SAS Institute Inc. 2021. SAS. https://www.sas.com/
- [25] Eshin Jolly. 2018. Pymer4: connecting R and Python for linear mixed modeling. Journal of Open Source Software 3, 31 (2018), 862.
- [26] Eunice Jun, Melissa Birchfield, Nicole De Moura, Jeffrey Heer, and Rene Just. 2022. Hypothesis Formalization: Empirical Findings, Software Limitations, and Design Implications. In ACM Transactions on Computer-Human Interaction (TOCHI), Vol. 29. Issue 1. "https://arxiv.org/abs/2104.02712"
- [27] Eunice Jun, Maureen Daum, Jared Roesch, Sarah E Chasins, Emery D Berger, Rene Just, and Katharina Reinecke. 2019. 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.
- [28] Moe Kayali, Babak Salimi, and Dan Suciu. 2020. Demonstration of inferring causality from relational databases with CaRL. Proceedings of the VLDB Endowment 13, 12 (2020), 2985–2988.
- [29] Ita GG Kreft, Ita Kreft, and Jan de Leeuw. 1998. Introducing multilevel modeling. Sage.
- [30] Erin LeDell and Sebastien Poirier. 2020. H2O AutoML: Scalable Automatic Machine Learning. 7th ICML Workshop on Automated Machine Learning (AutoML) (July 2020). https://www.automl.org/wp-content/uploads/2020/07/AutoML_ 2020_paper_61.pdf
- [31] Jiali Liu, Nadia Boukhelifa, and James R Eagan. 2019. Understanding the Role of Alternatives in Data Analysis Practices. IEEE transactions on visualization and computer graphics 26, 1 (2019), 66–76.

- [32] Yang Liu, Tim Althoff, and Jeffrey Heer. 2019. Paths Explored, Paths Omitted, Paths Obscured: Decision Points & Selective Reporting in End-to-End Data Analysis. arXiv preprint arXiv:1910.13602 (2019).
- [33] Yang Liu, Alex Kale, Tim Althoff, and Jeffrey Heer. 2020. Boba: Authoring and visualizing multiverse analyses. IEEE Transactions on Visualization and Computer Graphics (2020).
- [34] James Lloyd, David Duvenaud, Roger Grosse, Joshua Tenenbaum, and Zoubin Ghahramani. 2014. Automatic construction and natural-language description of nonparametric regression models. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 28.
- [35] Steson Lo and Sally Andrews. 2015. To transform or not to transform: Using generalized linear mixed models to analyse reaction time data. Frontiers in psychology 6 (2015), 1171.
- [36] Wendy E Mackay, Caroline Appert, Michel Beaudouin-Lafon, Olivier Chapuis, Yangzhou Du, Jean-Daniel Fekete, and Yves Guiard. 2007. Touchstone: exploratory design of experiments. In Proceedings of the SIGCHI conference on Human factors in computing systems. ACM, 1425–1434.
- [37] Ellicott C Matthay and M Maria Glymour. 2020. A graphical catalog of threats to validity: Linking social science with epidemiology. *Epidemiology (Cambridge, Mass.)* 31, 3 (2020), 376.
- [38] Richard McElreath. 2020. Statistical rethinking: A Bayesian course with examples in R and Stan. CRC press.
- [39] Muthén & Muthén. [n. d.]. MPlus. https://www.statmodel.com/
- [40] Brad A Myers, John F Pane, and Andy Ko. 2004. Natural programming languages and environments. Commun. ACM 47, 9 (2004), 47–52.
- [41] John Ashworth Nelder and Robert WM Wedderburn. 1972. Generalized linear models. Journal of the Royal Statistical Society: Series A (General) 135, 3 (1972), 370–384
- [42] Amadeu Anderlin Neto and Tayana Conte. 2013. A conceptual model to address threats to validity in controlled experiments. In Proceedings of the 17th International Conference on Evaluation and Assessment in Software Engineering. 82–85.
- [43] Don Norman. 2013. The design of everyday things: Revised and expanded edition. Basic books.
- [44] Judea Pearl. 1995. Causal diagrams for empirical research. Biometrika 82, 4 (1995), 669–688.
- [45] Judea Pearl. 1995. Causal diagrams for empirical research. Biometrika 82, 4 (1995), 669–688.
- [46] Judea Pearl et al. 2000. Models, reasoning and inference. Cambridge, UK: CambridgeUniversityPress 19 (2000).
- [47] Josef Perktold, Skipper Seabold, Jonathan Taylor, and statsmodels developers. 2020. Statsmodels v0.10.2 Reference Guide. (2020). "https://www.statsmodels.org/stable"
- [48] Adrian E Raftery. 1996. Approximate Bayes factors and accounting for model uncertainty in generalised linear models. *Biometrika* 83, 2 (1996), 251–266.
- [49] William S Robinson. 1950. Ecological correlations and the behavior of individuals. Sociological Review 15 (1950), 351–357.

- [50] Donald B Rubin. 2004. Teaching statistical inference for causal effects in experiments and observational studies. *Journal of Educational and Behavioral Statistics* 29, 3 (2004), 343–367.
- [51] Babak Salimi, Harsh Parikh, Moe Kayali, Lise Getoor, Sudeepa Roy, and Dan Suciu. 2020. Causal relational learning. In Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data. 241–256.
- [52] Michel F Sanner et al. 1999. Python: a programming language for software integration and development. J Mol Graph Model 17, 1 (1999), 57–61.
- [53] Skipper Seabold and Josef Perktold. 2010. Statsmodels: Econometric and statistical modeling with python. In Proceedings of the 9th Python in Science Conference, Vol. 57. Scipy, 61.
- [54] William R Shadish. 2010. Campbell and Rubin: A primer and comparison of their approaches to causal inference in field settings. Psychological methods 15, 1 (2010) 3
- [55] N.J.A. Sloane and R.H. Hardin. 2017. Gosset: A General-purpose program for designing experiments. http://neilsloane.com/gosset/
- [56] Peter Spirtes. 1994. Conditional independence in directed cyclic graphical models for feedback. (1994).
- [57] Peter Spirtes, Thomas Richardson, Christopher Meek, Richard Scheines, and Clark Glymour. 1996. Using d-separation to calculate zero partial correlations in linear models with correlated errors. *Publisher: Carnegie Mellon University* (1996)
- [58] IBM SPSS. 2021. SPSS Software. https://www.ibm.com/analytics/spss-statistics-software
- [59] Sara Steegen, Francis Tuerlinckx, Andrew Gelman, and Wolf Vanpaemel. 2016. Increasing transparency through a multiverse analysis. Perspectives on Psychological Science 11, 5 (2016), 702–712.
- [60] Etsuji Suzuki, Tomohiro Shinozaki, and Eiji Yamamoto. 2020. Causal diagrams: pitfalls and tips. Journal of epidemiology (2020), JE20190192.
- [61] Etsuji Suzuki and Tyler J VanderWeele. 2018. Mechanisms and uncertainty in randomized controlled trials: A commentary on Deaton and Cartwright. Social science & medicine (1982) 210 (2018), 83–85.
- [62] Emi Tanaka. 2021. Edibble: An R-package to construct designs using the grammar of experimental design. https://github.com/emitanaka/edibble
- [63] R Core Team et al. 2013. R: A language and environment for statistical computing. (2013).
- [64] Johannes Textor, Juliane Hardt, and Sven Knüppel. 2011. DAGitty: a graphical tool for analyzing causal diagrams. *Epidemiology* 22, 5 (2011), 745.
- [65] C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown. 2013. Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms. In Proc. of KDD-2013. 847–855.
- [66] Tyler J VanderWeele. 2019. Principles of confounder selection. European journal of epidemiology 34, 3 (2019), 211–219.
- [67] Priscilla Velentgas, Nancy A Dreyer, Parivash Nourjah, Scott R Smith, Marion M Torchia, et al. 2013. Developing a protocol for observational comparative effectiveness research: a user's guide. (2013).