

# Measurement as Bricolage: Examining How Data Scientists Construct Target Variables for Predictive Modeling Tasks

LUKE GUERDAN<sup>†</sup>, Carnegie Mellon University, USA

DEVANSH SAXENA<sup>†</sup>, Carnegie Mellon University, USA

STEVIE CHANCELLOR\*, University of Minnesota, USA

ZHIWEI STEVEN WU\*, Carnegie Mellon University, USA

KENNETH HOLSTEIN\*, Carnegie Mellon University, USA

Data scientists often formulate predictive modeling tasks involving fuzzy, hard-to-define concepts, such as the “authenticity” of student writing or the “healthcare need” of a patient. Yet the process by which data scientists translate fuzzy concepts into a concrete, proxy target variable remains poorly understood. We interview fifteen data scientists in education ( $N=8$ ) and healthcare ( $N=7$ ) to understand how they construct target variables for predictive modeling tasks. Our findings suggest that data scientists construct target variables through a *bricolage* process, involving iterative negotiation between high-level measurement objectives and low-level practical constraints. Data scientists attempt to satisfy five major criteria for a target variable through bricolage: validity, simplicity, predictability, portability, and resource requirements. To achieve this, data scientists adaptively use *problem (re)formulation strategies*, such as *swapping* out one candidate target variable for another when the first fails to meet certain criteria (e.g., predictability), or *composing* multiple outcomes into a single target variable to capture a more holistic set of modeling objectives. Based on our findings, we present opportunities for future HCI, CSCW, and ML research to better support the art and science of target variable construction.

**CCS Concepts:** • Human-centered computing → Empirical studies in HCI; • Computing methodologies → Machine learning; Supervised learning.

**Additional Key Words and Phrases:** data science, interview study, measurement, validity, model evaluation

## 1 INTRODUCTION

*“Problem-solvers who do not proceed from top-down design but are arranging and rearranging a set of well-known materials can be said to be practicing bricolage. They tend to try one thing, step back, reconsider, and try another”* (Turkle, p. 51).

Data scientists grapple with subtle but important considerations when translating organizational goals into a formal modeling task. For example, consider a data science team developing a model to help physicians decide which patients have sufficient “healthcare need” for enrollment in a high-risk care management program (e.g., [102]). While developing a model, data scientists may work with clinicians to understand how to translate risk predictions to actionable interventions – for example, by selecting an appropriate risk cut-off for enrollment recommendations. Data scientists may also carefully consider whether the patient populations represented in their training data adequately match the characteristics of patients encountered when the model is deployed. The process of *problem formulation* describes how data scientists navigate such translational concerns as they map a high-level goal into a tractable computational problem [106].

One critical aspect of problem formulation is the process of *target variable construction* (Figure 1). Many outcomes of interest to data scientists – such as the “authenticity” of student writing or the “healthcare need” of patients – are latent, theoretical constructs that are not directly observable in

---

<sup>†</sup> Co-first authors contributed equally to this research. \*Co-senior authors contributed equally to this research.

Authors’ Contact Information: Luke Guerdan<sup>†</sup>, lguerdan@cs.cmu.edu, Carnegie Mellon University, Pittsburgh, PA, USA; Devansh Saxena<sup>†</sup>, Carnegie Mellon University, Pittsburgh, PA, USA; Stevie Chancellor\*, University of Minnesota, Minneapolis, MN, USA; Zhiwei Steven Wu\*, Carnegie Mellon University, Pittsburgh, PA, USA; Kenneth Holstein\*, Carnegie Mellon University, Pittsburgh, PA, USA.

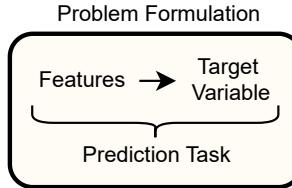


Fig. 1. An illustration of the relationship between the target variable and problem formulation. A *prediction task* describes a mapping from *features* to a *target variable*. A *problem formulation* connects a prediction task to high-level modeling goals, including factors such as how risk scores are converted to intervention recommendations and the target population at deployment time.

data. To develop a model targeting these constructs, data scientists must construct a *proxy variable* based on outcomes that are readily available. For example, a data scientist might choose to predict patients' cost of medical care or their number of recently diagnosed chronic health conditions as a proxy for their "healthcare need" [102]. Traditionally, measurement theory assumes that analysts follow a top-down workflow while mapping unobserved constructs to measurements in data [62]. For example, a learning scientist might identify a latent construct they would like to measure among students (e.g., "math proficiency") and design an operationalization (e.g., a test) *before* collecting data (e.g., student test scores). Disciplines in the quantitative social sciences have developed a rich set of guidelines (e.g., [90, 138]) and modeling tools (e.g., [38, 129]) to help analysts implement this top-down measurement workflow.

**However, the traditional top-down conceptualization of measurement is in tension with our understanding of data science as a bottom-up process constrained by existing data** [77, 106]. Passi and Barocas [106] observe that problem formulation in data science is an iterative process, where data scientists progressively reconcile their high-level objectives with low-level data and modeling constraints. As a result, while a traditional top-down measurement approach assumes that an analyst has the resources to collect data tailor-made to the goals of a study, data scientists are often forced to re-purpose the limited data that happens to be available. This disconnect creates a significant challenge: though many flawed data science projects can be traced back to target variable construction (e.g., [7, 26, 72, 94, 136]), existing tools, frameworks, and pedagogical resources reflect a top-down planning approach towards measurement, rarely accounting for the bottom-up, resource-constrained reality of data science practice.

To understand data scientists' current practices and challenges in constructing target variables—and to identify opportunities to better support target variable construction in practice—in this work, we conducted semi-structured interviews with fifteen data scientists from education (N=8) and healthcare (N=7). We asked data scientists to walk us through how they navigate target variable construction in prior modeling projects. Participants also engaged in a target variable construction task, allowing us to observe their reasoning in a new modeling context. Our goal in these interviews was to understand how data scientists intertwine two conceptually distinct activities, measurement of latent constructs and prediction of target variables, during problem formulation [97]. While prior empirical studies of data science practice have often centered prediction (e.g., [6, 69, 143]), we broaden our scope of inquiry to investigate the interplay between measurement validity and predictive utility in problem formulation.

We find that data scientists construct target variables by **leveraging creative and pragmatic approaches to make do with the limited data they have at hand**. Data scientists craft their formulation through a process of iterative *dialogue* with their data, tweaking and refining their

target variable until they conclude that it provides a sufficient operationalization for the unobserved construct of interest. To understand data scientists' target variable construction practices, we draw upon the concept of *bricolage*, originally introduced by Claude Lévi-Strauss to contrast the rigid top-down methodology of Western science with the more tangible "science of the concrete" practiced in non-Western societies [81]. Rather than implementing abstract concepts through top-down design, the bricoleur solves problems by arranging and re-arranging well-known materials in a bottom-up process. Like bricoleurs, data scientists in our study often acknowledged imperfection in their improvisational approach to measurement, yet they also identified innovative solutions that would be difficult to anticipate through a process of top-down design.

Based on our findings, **we synthesize a new process model that better reflects data scientists' bricolage approach towards measurement** (Figure 2). Data scientists balance the validity of their formulation with other criteria such as predictive performance, simplicity, portability, and resource requirements. Data scientists narrow in on a suitable formulation by drawing upon a rich repertoire of problem (re)formulation strategies. Data scientists apply a strategy, take a step back to assess their progress along each criterion, and then apply a new strategy until all evaluation criteria are satisfied. Data scientists sometimes discontinue their projects if they conclude that no strategies are sufficient given the resources at hand.

Our characterization of target variable construction as bricolage **offers an important counterbalance to nascent views on measurement of AI systems**. While recent work correctly identifies measurement problems as the source of many AI failures [5, 28, 47, 62, 71, 83, 114, 115], proposed solutions often assume or reinforce a top-down, idealized measurement process. Our findings suggest that in order for established measurement principles to have impact in practice, there is a need for thoughtful reconciliation with the resource-constrained, bottom-up reality of data science work. By documenting how measurement and prediction intertwine in data scientists' bricolage practices, we contribute a more grounded understanding of how measurement actually occurs in the wild, providing a foundation for developing more effective interventions. While discussing implications of our findings, **we present new opportunities for future CSCW, HCI, and ML research to develop tools that better support data scientists' bricolage practices**.

In summary, our contributions to the CSCW literature are:

- A novel characterization of measurement in data science as a *bricolage practice*. We illustrate how data scientists adaptively balance multiple evaluation criteria by creatively combining existing data under resource constraints.
- Empirical findings from interviews with fifteen data scientists in healthcare and education, synthesized into a process model that reveals how practitioners weigh competing criteria (§ 4.1), apply five strategic approaches to (re)formulation (§ 4.2), and evaluate validity through both theory-driven and data-driven practices (§ 4.3).
- Design implications for how the CSCW, HCI, and ML research communities might better scaffold data scientists' bricolage practices going forward. These recommendations explore how we might enhance the validity of data scientists' resulting formulations while also supporting the flexibility that makes bricolage effective (§ 5).

## 2 BACKGROUND AND RELATED WORK

In this section, we begin by introducing Claude Lévi-Strauss' theory of Bricolage. We then describe existing Science and Technology Studies (STS), HCI, and CSCW research that has investigated questions related to categorization, measurement, and target variable construction in data science practice. We conclude by providing an overview of key concepts from measurement theory in the quantitative social sciences which we build upon in this work.

## 2.1 Bricolage and the Science of the Concrete

*Bricolage*, defined by Claude Lévi-Strauss as “making the most of available resources”[81], describes how actors (i.e., *bricoleurs*) construct artifacts in creative and adaptive ways. Bricolage is central to what Lévi-Strauss calls *the science of the concrete*, which contrasts with modern scientific thought by focusing on the materials directly at hand rather than abstract principles. Bricolage contains three elements: repertoire, dialogue, outcome [36]. *Repertoire* describes the resources—physical artifacts, ideas, or skills—that the bricoleur brings to bear on a problem. Whereas Lévi-Strauss viewed the traditional engineer or scientist as acquiring resources tailor-made to specification, the bricoleur slowly assembles various “odds and ends” that might be useful in the future. Then, when faced with a new problem, the bricoleur engages in *dialogue* with these materials—a process of arranging and re-arranging materials in her repertoire to find a workable solution. While engaging in dialogue, the bricoleur adopts old objects for new purposes and combines concepts with materials in novel ways. Finally, *outcome* describes the artifact resulting from the bricolage process. While bricolage entails a tolerance for imperfection, Lévi-Strauss argues that bricoleurs sometimes identify “brilliant unforeseen results” [81] as part of their process.

Prior research has used Lévi-Strauss’ framework to describe many forms of knowledge work as bricolage (e.g., entrepreneurship [36, 88], design [19, 84, 137], management [36], and programming [135]). For instance, Baker and Nelson [9] use bricolage to explain how entrepreneurs build business ventures under resource constraints. They observe that entrepreneurial success hinges on the ability of firms to combine their limited resources in creative, adaptive ways.

In her book *Life on the Screen*, Sherry Turkle uses bricolage to describe how programmers write computer code. Through the mid-1980s, the canonical coding approach taught in classrooms was *structured programming*—a linear, top-down process that proceeds by systematically mapping abstract plans into procedural instructions. Yet, Turkle observed that programmers—novices and experts alike—often engaged in a more unstructured *dialogue* with their code while solving real-world problems. Drawing an analogy between the programmer and bricoleur, she explains:

*“In the context of programming, the bricoleur’s work is marked by a desire to play with lines of computer code, to move them around almost as if they were material things—notes on a score, elements of a collage, words on a page”* (Turkle, p. 51-52)

Turkle describes this style of tinkering as important for developing an intuition for the complexity of computer programs. Like a painter who steps back to inspect a canvas between brush strokes, programmers try one approach, re-evaluate, and then try another. Turkle observes that this process often leads programmers to more elegant solutions than could be imagined through the more traditional process of structured planning.

In this work, we draw a novel connection between *data science work* and the practice of bricolage. Much like an entrepreneur, programmer, or painter, data scientists construct target variables by combining *existing* resources in creative, adaptive ways. Instead of instantiating abstract plans through top-down design—akin to the structured programming approach taught in classrooms—data scientists engage in more direct *dialogue* with their data and modeling tools, making do with the materials at hand. And analogous to past research investigating how to support such exploratory and improvisational coding practices [73], this has implications for the design of practical tools to support target variable construction in practice. We highlight connections between target variable construction and bricolage throughout the findings (§ 4) and describe implications for tool design in the discussion (§ 5).

## 2.2 Tracing the Social Construction of Target Variables

Target variable construction also relates to a rich body of literature in Science and Technology Studies (STS) and philosophy of science that investigates the social processes that give rise to categories in data [16, 29, 49, 52, 79, 93, 110, 111, 132]. Latour and Woolgar [79] provide an account of the social construction of scientific facts within laboratories, highlighting that facts are not merely discovered but actively produced through a complex interplay of social and technical factors. Likewise, Bowker [16] examines how systems of classification, such as the International Classification of Diseases (ICD), shape and are shaped by the social and political contexts in which they are embedded. Building on Bowker’s work, recent studies have also explored the construction of racial and ethnic categories within AI models (e.g., [2, 89]) and charted measurement assumptions underpinning motion capture technologies [52].

Most related to our study, Muller et al. [96] interviewed fifteen data scientists to investigate the *data labeling* and *annotation* practices of data science teams. Muller et al. [96] identified three approaches that teams used while defining labels: *principled*, *iterative*, and *improvisational*. Teams engaging in principled design followed a top-down planning approach, systematically specifying annotation guidelines before data labeling. Teams leveraging iterative design took a bottom-up approach, incorporating opportunities for annotation protocol refinement throughout the labeling process. Finally, teams engaging in improvisational design viewed data labeling as an open-ended, exploitative process. Our study investigates a distinct question of how data science teams construct target variables for *predictive modeling tasks*. Rather than studying how data science teams design “ground truth” labels anew, we examine how data scientists connect their *existing* data to their high-level modeling objectives as part of the problem formulation process.

## 2.3 Empirical Studies of Data Science Practice

Our work also builds upon empirical studies investigating how data scientists conduct their day-to-day work. Many studies have examined activities in the technical “inner loop” [77] of data science practice, such as exploratory data analysis (e.g., [6, 143]), feature engineering (e.g., [95]), and visualization (e.g., [69]). Other studies have investigated broader “outer loop” activities performed by data scientists, including collaborating with other stakeholders [60, 75, 77, 87, 99, 107, 147] and addressing ethical concerns [31, 59]. CSCW research has described inner and outer loop activities as iterative and interconnected [77, 106]. In particular, based on interviews with ten professionals spanning industry and academia, Kross and Guo [77] observed that data scientists iteratively “framed the problem” (i.e., an outer loop activity) and returned to their data to gain an understanding of its affordances (i.e., inner loop activities).

However, we currently have an incomplete understanding of *how* data scientists construct target variables while formulating prediction tasks. To date, Passi and Barocas [106] provide the most detailed account of how problem formulation unfolds in real-world modeling contexts. Over the course of an eighteen-month ethnographic study of a corporate data science team, Passi and Barocas [106] observed that data scientists repeatedly (re)formulated their prediction task as they discovered new information about data availability constraints and the predictability of alternative outcome variables. Based on their findings, Passi and Barocas [106] describe problem formulation as a *negotiated translation* between high-level objectives and a concrete prediction task. Our work builds upon this finding to identify the more general evaluation desiderata and (re)formulation strategies driving the problem formulation process. We develop a process model through analysis of our findings and use this model to identify targeted opportunities to better support data scientists’ problem formulation practices going forward. For example, when viewed through the lens of our findings, the team studied by Passi and Barocas [106] used *one of five*

(re)formulation strategies (i.e., *swapping*) to change outcome variables when the first proved infeasible due to data availability constraints. We present design opportunities to help data scientists more effectively apply swapping (§ 4.2.3) and other (re)formulation strategies.

## 2.4 Toolkits Supporting Target Variable Construction

The limitations of data science toolkits and analysis methods have been documented across empirical studies examining data science practices [64], software repositories [57], and ML-based sciences [70]. Jun et al. reviewed 20 commonly-used data analysis tools and found that the ecosystem of tools supports vastly different model implementations, even when using the same underlying mathematical equation [64]. Although these tools support nuanced model implementations, their low-level technical abstractions make it challenging for analysts to engage in the high-level conceptual reasoning required for hypothesis formalization [64]. As a result, interpretations of data, modeling assumptions, and implementation details can significantly influence outcomes and impact the validity of results [57, 64, 65, 70]. Research examining software engineering practices in code repositories, such as GitHub, further demonstrates how the choice of analysis tools can lead to substantial differences in simple metrics (e.g., commit counts and developer contributions) [51, 57, 80]. This further poses threats to validity of findings derived from one analysis tool. Given this gap between technical abstractions and conceptual reasoning—and the embedded assumptions inherent in machine learning methods like feature engineering—it is not surprising that ML-based sciences continue to face a reproducibility crisis [70].

To mitigate these issues, the HCI, ML, and FAccT research communities have developed a range of toolkits to intervene at points along data scientists' inner and outer loop workflows. For instance, researchers have developed systems facilitating data collection and configuration (e.g., [13, 48, 68, 116]), exploratory data analysis (e.g., [74, 144]), model training (e.g., [33, 105]), caching intermediate states during data science experiments (e.g., [123]), and performance assessment (e.g., [21, 146]). Many tools have also been developed to help data scientists explain their models (e.g., [85, 117, 118]) or assess fairness-related concerns (e.g., [4, 12, 20, 33, 82, 119, 134, 141, 142]). In the FAccT community, researchers have also developed responsible AI checklists and guidelines to help data scientists assess and document components of their AI system (e.g., dataset information [42, 112], intended use cases [92]). To our knowledge, Gala et al. [41] propose the only system specifically designed to support the target variable construction process. Instead, **existing toolkits are often designed to support data scientists' practices only after a target variable has been specified**. This presents a pressing gap in data science support at a critical point in their workflows.

However, prior HCI research suggests that such toolkits may prove ineffective if they are not designed around the *existing* practices of data scientists. For instance, past efforts to improve programming tools without a human-centered assessment of programmers' needs yielded ineffective languages and systems [54, 76, 98, 125]. Similarly, the utility of ML fairness toolkits was found to be limited in practice as they were designed for isolated use by technical roles, while actual fairness work often relied on deep cross-functional collaborations [31]. Our work aims to help the research community avoid similar challenges while developing tools to support data scientists' problem formulation processes. In particular, we identify opportunities for the HCI, CSCW, and ML research communities to develop tools that help data scientists more carefully weigh trade-offs across multiple evaluation criteria (§ 5.3.1), select (§ 5.3.2) and apply (§ 5.3.3) strategies, and evaluate the validity of their problem formulation more effectively (§ 5.3.4).

## 2.5 Measurement and Validity in AI

A growing line of research traces AI system failures to *validity* pitfalls [28, 45, 47, 62, 91, 94, 102, 131]. In the quantitative social sciences, the validity of a measurement instrument is established via a multifaceted assessment along several dimensions. *Construct validity* describes the extent to which a measurement instrument wholly and fully reflects the theoretical construct it purports to measure [62]. In predictive modeling, construct validity is most often discussed in connection to whether the target variable predicted by a model reflects a theoretical construct of interest to model developers [28, 47].<sup>1</sup> *Internal* and *external* validity are also important dimensions in experimental research. *Internal validity* describes the extent to which an experimental study can confidently establish a causal relationship between independent and dependent variables [43]. In predictive modeling, it relates to the existence of a defensible causal relationship between features and a target variable [28]. *External validity*, on the other hand, measures how well an analysis can be generalized beyond the specific conditions of the study. This dimension has been linked to the extent to which a predictive model’s in-distribution performance generalizes across real-world deployment contexts [28].

While these validity dimensions articulate helpful conceptual desiderata for a predictive model, they provide a limited picture of *how* practitioners actually navigate measurement decisions in their day-to-day work. Mussgnug [97] argue that current practice is characterized by a “predictive reframing”, in which data scientists cast a measurement problem (inferring existing but unobserved quantities) as a prediction task. This reframing shifts the epistemic aim from measuring a concept to predicting a given measurement of that concept. For example, when data scientists attempt to assess poverty levels using satellite imagery, they often frame their task as “predicting poverty” rather than “measuring poverty.” This practice centers the evaluation of a machine learning model around narrow predictive performance characteristics while sidestepping deeper questions of validity.

We expand upon Mussgnug [97]’s argument by characterizing how data scientists weigh validity and predictive performance during problem formulation. In support of Mussgnug [97]’s argument, we find that many prediction tasks pursued by data scientists are indeed measurement problems. For example, some participants in our study used supervised learning models to measure concepts recorded at future points in time – e.g., “10-year cardiovascular risk.” Others used supervised learning models as measurement instruments while measuring a concept at a current point of time – e.g., to score the “authenticity” of a student’s writing. In both cases, data scientists heavily anchored on predictive performance while evaluating their problem formulation. However, counter to Mussgnug [97]’s argument, we find that data scientists *do* also engage with validity concepts during problem formulation, albeit with a more pragmatic, situated approach that lacks formal terminology. We uncover how data scientists evaluate validity in our findings (§ 4), and unpack implications for data science pedagogy (§ 5.2.3) and tooling support (§ 5.3.4) in the discussion.

## 3 METHODS

To examine how data scientists construct target variables for predictive modeling tasks, we interviewed 15 data scientists from the education (N=8) and healthcare (N=7) sectors who had prior experience with developing and evaluating predictive models. We conducted one-hour semi-structured interviews to understand participants’ current practices and perceived challenges for measurement as they engaged in problem formulation.

---

<sup>1</sup>Jacobs and Wallach [62] provide a detailed review of sub-dimensions of construct validity as they pertain to measuring fairness properties of algorithmic systems.

PID	Domain	Postgrad Experience	Current Primary Role	Current Affiliation
P1	Education	15 years	Board Member	Industry
P2		7 years	Senior Analyst, Ph.D. Student	Industry/Academic
P3		5 years	Postdoc	Academic
P4		2 years	Ph.D. Student	Academic
P5		2 years	Ph.D. Student	Academic
P6		11 years	VP of Data Science	Industry
P7		10 years	Data Scientist	Industry
P8		11 years	Senior Research Scientist	Industry
P9	Healthcare	7 years	Senior Research Scientist	Industry
P10		6 years	Research Associate	Academic
P11		4 years	Senior Data Scientist	Industry
P12		11 years	Senior Data Science Analyst	Industry
P13		8 years	Senior Data Scientist	Industry
P14		8 years	PostDoc	Academic
P15		10 years	Assistant Professor	Academic

Table 1. Demographic overview of study participants: data scientists' experience, roles, and affiliations.

### 3.1 Study Design

3.1.1 *Directed Storytelling.* During the first interview segment, we used a directed storytelling approach [37], inviting data scientists to reflect on a time in the past when they developed and evaluated a predictive model. We centered questions in this segment around data scientists' *target variable construction process*, a component of problem formulation which prior literature has shown is both complex and consequential [102].

To ground our discussions, we first asked participants to share a few concrete experiences where they had developed and/or evaluated predictive models as part of their role. Building on these concrete experiences, the interviewer then narrowed in on one of the projects mentioned and asked participants the following follow-up questions, in a semi-structured fashion:

- (1) Were there ever times when you felt uncertain about your choice of target variable?
- (2) Could you walk me through how your team went about picking a target variable?
- (3) How did you know whether your team selected the “right” target variable?

We also asked participants a range of additional follow-up questions when these were not already covered by participants during their earlier responses, such as measurement guidelines and best practices they consulted throughout their project. On average, the first segment lasted 30-40 minutes of the interview.

3.1.2 *Vignette-Based Problem Formulation Task.* During the second segment, we asked participants to reason about a hypothetical problem formulation as we walked them through a vignette. We aimed to better understand how participants reason about measurement decisions on-the-fly, while potentially surfacing challenges and opportunities they might not have recalled while reflecting on their prior experiences. We asked participants with a healthcare background to imagine they were developing a model to predict patients with high future “medical acuity,” and we asked participants with an education background to imagine they were developing a model to predict which students

were “academically at-risk” (see Appendix A).<sup>2</sup> We provided participants a set of features and several outcome variables and instructed them to imagine that these were the data available for model development.

During the vignette, we asked participants to reflect on which outcomes might serve as a “better or worse choice” of prediction target given the stated modeling goals. After learning how participants initially approached this task, we then presented a series of model evaluation plots indicating the performance of predictive models targeting each outcome (see Appendix A). After showing each plot, we paused and asked participants to again reflect on which outcomes might serve as better or worse choices for the stated modeling goals. We also asked participants to share any additional information that they would ideally like to have, to better inform their thought process. As participants engaged with the task over subsequent interviews, we refined the specificity of the scenario description and accompanying evaluation plots to further probe issues that surfaced in prior interviews. This second segment lasted 20-30 minutes on average.

### 3.2 Participants

Participants were recruited via purposive sampling techniques [56] from authors’ professional networks and social media (Table 1). To ensure that participants had a baseline level of domain expertise, we required that data scientists had domain-specific graduate-level training related to healthcare or education (e.g., biostatistics, learning sciences) or at least one year of full-time industry experience in a relevant domain to be eligible to participate. We further required that participants had prior experience developing and evaluating predictive models as part of their role. Our participants had a range of post-graduate experience at the time of the interview, spanning from 3-6 months to 15 years. At the time of the interview, participants had a variety of primary roles (e.g., Data Scientist, Assistant Professor, Board Member, PostDoc) from both academic and industry affiliations. As a result, our sample consists of data scientists with a baseline level of domain knowledge and a range of technical experience. This study was approved by the University Institutional Review Board and all participants were compensated with a \$35 Amazon gift card for completing the study.

### 3.3 Data and Analysis

One of the authors conducted interviews remotely via Zoom between September and December 2023. In total, we recorded and transcribed 14 hours and 30 minutes of audio after collecting verbal and written consent from all participants. Two authors independently performed open coding on transcripts and met regularly to discuss and reconcile interpretation differences. After interviewing fifteen participants, we observed a common set of themes recurring across interviews and took this as an indicator that we had reached saturation [120]. We adopted a reflexive thematic analysis approach [18, 27] to group codes into successively higher themes, informed by discussions among all of the authors. Our process gave rise to sixteen high-level themes, organized into three categories: (1) *factors* driving the target variable construction process, which we discuss in §4.1; (2) data scientists’ *model development* and *evaluation* practices, which we discuss in §4.2 and §4.3; and (3) data scientists’ *orientation* towards the problem formulation process, which we discuss throughout the findings, when drawing connections to the concept of bricolage.

---

<sup>2</sup>P1 was shown the healthcare vignette due to data availability constraints encountered at early phases of the interview schedule. During the interview, P1 signaled places where their interpretation of the vignette was informed by their domain experience. We took this into consideration during analysis.

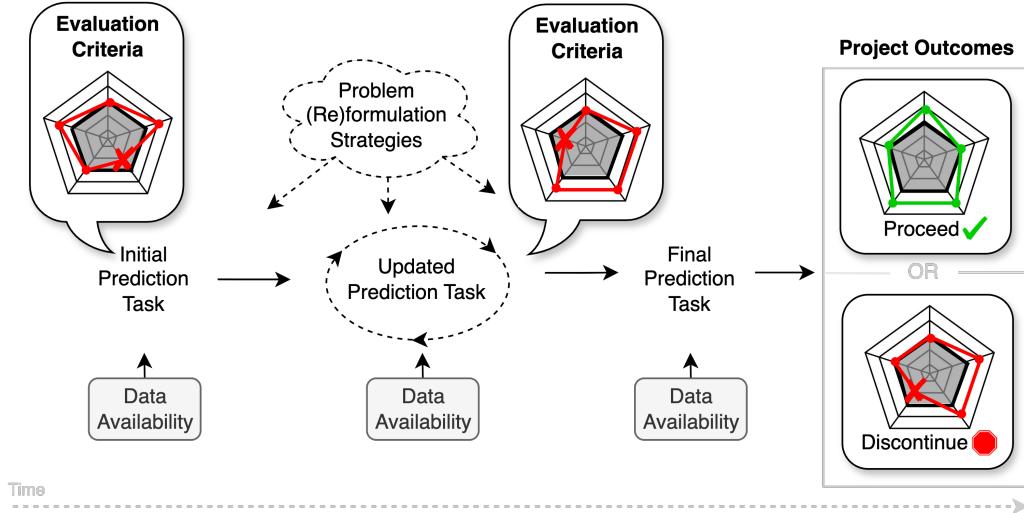


Fig. 2. An illustration of the target variable construction process presented in our findings. Data scientists specify an initial prediction task based on their available data, then iteratively refine their prediction task by applying (re)formulation strategies. Data scientists proceed with their final prediction task if it satisfies all criteria, or discontinue their project if strategies are exhausted.

### 3.4 Positionality

Our positionality as a research team is shaped by our backgrounds spanning human-computer interaction, cognitive psychology, machine learning, and design. Several team members have professional experience as data scientists, while others have training in educational analytics and behavioral health. This interdisciplinary background helped us recruit participants and build rapport during interviews. During analysis, our team's quantitative social sciences background helped us connect the practices we documented to established ideas in validity theory. Our interdisciplinary perspective also directed our focus toward tensions between pragmatic data science realities—which we personally encountered—and formal validity concepts from academic research. Team members with design experience had previously encountered Claude Lévi-Strauss' concept of bricolage, which became our central theoretical framework. The pragmatic orientation in our findings emerged from combining academic social science training with real-world data science experience.

## 4 FINDINGS

Our findings illustrate that target variable construction is a bricolage practice, in which data scientists design measurement instruments (i.e., target variables) for unobservable constructs by creatively making do with their existing data. While measurement has traditionally been conceptualized as a *monologue* [135] – a one-way, top-down translation from theoretical constructs to measurements in data – we find that data scientists instead construct target variables through a rich process of *dialogue* with their data, driven by a need to make do with resource constraints.

The specific process of bricolage unfolds as depicted in Figure 2. Data scientists begin by formulating an initial prediction task based on available data. They then iteratively inspect their formulation along several dimensions: “Can the target variable be predicted with ‘reasonable’ accuracy?”, “Will the necessary data be available across deployment contexts?” When issues arise along one or more criteria, data scientists apply various (re)formulation strategies to address these

gaps. These strategies serve as the “brushstrokes” of data science work. Like a painter with a canvas, data scientists apply a strategy, step back to survey their progress along multiple criteria, then consider a new strategy in response to revealed flaws. The resulting problem formulation, while rarely ideal, reflects what the bricoleur considers “fit for purpose”, by acknowledging imperfection while identifying creative and practical solutions that might never emerge through rigid top-down design. Data scientists eventually proceed with a formulation if they believe it satisfies their criteria across multiple dimensions, or discontinue their project if all strategies have been exhausted.

In the following subsections, we detail this bricolage process by describing how data scientists balance validity with other competing criteria (§ 4.1), the (re)formulation strategies they apply to navigate this balancing act (§ 4.2), and the theory-driven and data-driven approaches they use to evaluate the validity of their formulation (§ 4.3).

## 4.1 Balancing Validity With Other Criteria

While engaging in target variable construction, data scientists were forced to make do with the data they already had available for the task at hand. This process of making do forced data scientists to navigate a rich space of tradeoffs between validity and other important criteria, such as predictive performance, simplicity, resource requirements, and the ease by which it could be translated to different geographic or institutional contexts (i.e., portability) (Figure 3). In Table 2, we provide an overview of illustrative tradeoffs that data scientists made between validity and other criteria. We provide a more detailed description of how the process of making tradeoffs unfolded in the sections below.

**4.1.1 Balancing Validity with Simplicity.** While constructing prediction tasks, several data scientists navigated a tension between the *validity* and *simplicity* of their formulation. In particular, data scientists found complex formulations, such as those combining multiple outcome variables, hard to explain to other stakeholders and maintain in production. For example, P1 was a director of a data science team at a large education company. Because the essay scoring tool designed by their team was used by schools throughout the country, P1 described the *simplicity* of their formulation as critical for ensuring the system was robust and maintainable. P1 *had* considered adopting a sophisticated multi-target formulation to improve validity. Yet, their team later abandoned the approach over practical concerns, explaining that the sophisticated models “*are a real hassle in practice. The extra [...] bells and whistles [...] made it flakier code, something that we didn’t necessarily want to deploy at scale, compared to a relatively straightforward single output variable model.*” Further, P1 later shared that they felt it was possible to proceed with using a simple single output variable model if they made several hyperparameter adjustments and user interface tweaks. In this experience, P1 was willing to settle for a specification that had minor validity drawbacks because they felt it would be simple and easy to maintain at scale.

**4.1.2 Balancing Validity with Resource Requirements.** Data scientists also evaluated the validity of their formulation against its resource requirements. For example, P6 was a data scientist developing an online math tutoring platform. P6 explained that the time and cost involved with acquiring “gold standard” student test scores to evaluate their models was a frequent challenge. School districts often had data sharing restrictions. Opportunities to test students were limited. However, P6 explained that even if they *could* test students more often, they would not recommend that approach. Testing students is “*just removing opportunity. Removing time from the students’ school calendar, which is already [...] crowded.*” Instead of collecting data, P6 devised a creative workaround that achieved their project’s goals using the data they *already* had available.<sup>3</sup> While P6 acknowledged that their

<sup>3</sup>We detail the strategy P6 and their team used in Section 4.2.4.

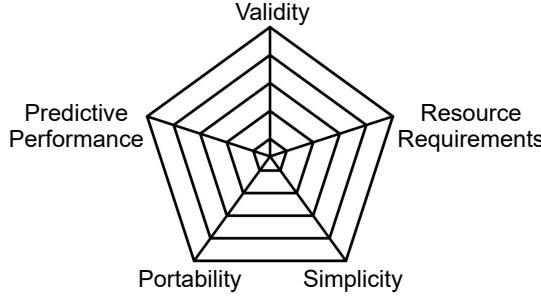


Fig. 3. Data scientists evaluated the validity, predictive performance, portability, simplicity, and resource requirements of their problem formulation.

Criterion	Description	Illustrative Tradeoff
<b>Validity</b>	The extent to which a model assesses the specific concept or construct that the data scientist intends to measure [28, 34].	N/A
<b>Predictive Performance</b>	The extent to which a model's predictions align with observed outcomes (e.g., accuracy or AU-ROC).	Predicting cardiovascular rather than renal diabetes health outcomes (P10). See § 4.2.3.
<b>Resource Requirements</b>	The organizational resources required to implement and maintain the formulation.	Making do with limited test scores to avoid imposing additional testing (P6). See § 4.1.2.
<b>Portability</b>	The ease by which a task definition can be transferred across geographic or institutional contexts.	Using coarse veterinary health information to support model inference across clinical contexts (P13).
<b>Simplicity</b>	The complexity of a task definition or model's structure and interpretation.	Switching from a complex multi-outcome model due to challenges in explanation (P2). See § 4.2.2.

Table 2. A description of criteria that data scientists considered while evaluating their problem formulation and illustrative examples of how participants balanced each with validity.

workaround was imperfect, they were fine with proceeding so long as they had “*at least some confidence*” in its validity. In particular, while reflecting on potential drawbacks of their modeling decisions, P6 shared that “*this [application] isn’t high stakes*” when compared to other potential modeling contexts (e.g., university admissions). However, had the model been deployed in other contexts, they explained that they might apply a different level of stringency to the validity of their approach.

**4.1.3 Elasticity of Evaluation Criteria.** We observed that data scientists applied varying levels of flexibility to different evaluation criteria, a characteristic we term *elasticity*. Participants frequently demonstrated elasticity with validity standards, tolerating imperfections when doing so reduced the complexity or resource requirements of their formulation. Data scientists actively *calibrated* this tolerance based on the success of compensatory adjustments (e.g., user interface tweaks) and the perceived “stakes” of their modeling decisions. Interestingly, participants maintained more rigid standards for other criteria. For instance, several participants cited hard cutoffs for acceptable predictive performance — e.g., P8 applied an AU-ROC threshold of 0.7 to determine whether to proceed with a project. This rigidity in predictive performance evaluation likely stems from the ease with which it can be distilled into a single, quantifiable metric, unlike validity, which requires a more holistic assessment involving multiple sources of evidence [62]. While tradeoffs are an inherent component of bricolage, motivated by adaptive response to resource constraints, the process of *navigating* tradeoffs remains poorly supported. We describe opportunities to help data scientists navigate tradeoffs in Section 5.3.1.

## 4.2 Strategies for (Re)formulating Prediction Tasks

Given their limited and imperfect data resources, how do data scientists engage in bricolage to craft a formulation that satisfies their objectives? Much like a painter refining a painting through successive brushstrokes, or a programmer compiling and re-compiling a system after adjusting lines of code [135], data scientists iteratively construct target variables through *(re)formulation strategies*. These strategies form the backbone of the dialogue between the data scientist and their formulation, enabling them to respond adaptively to constraints encountered while working with existing data. At times, data scientists apply these strategies instinctively, without a second thought. At other times, they shared eureka moments they encountered after discovering a strategy that leveraged unexpected data sources or combined information in novel ways. In the following sections, we detail five key strategies—piggybacking, composing, swapping, bridging, and refining—that data scientists used to navigate the constraints of their limited repertoire.

**4.2.1 Piggybacking.** During the initial stages of the problem formulation process, data scientists often constructed target variables through **piggybacking**: *adopting or extending a problem formulation that was previously used in a similar context in the past*. Participants justified piggybacking by citing precedents from related projects. For example, while developing a model to inform university admissions decisions, P5 suggested predicting whether students’ college GPA exceeded a 3.0 because other data science teams used it for similar admissions pipelines in the past. Piggybacking provided a rationale for P5’s formulation decisions and offered a point of comparison with prior work. However, they were also wary that piggybacking could introduce unintended validity consequences.

One participant encountered a validity flaw in their formulation after piggybacking. P1 was developing an automated essay scoring tool on a six-point scale, where six is the best. P1 initially piggybacked by adopting the most prominent evaluation metric in the literature (Quadratic Weighted Kappa) but quickly realized it led to jarring results.

*“The model simply never gave that top score. The more accurate [model] predicted the two’s, three’s, four’s, five’s. There are students out there that simply don’t accept that. They’ll write 30, 40, 50, 60 attempts at an essay, trying to get [...] a perfect score.”*

Piggybacking off of other modeling contexts did not work for P1’s dataset - optimizing for Quadratic Weighted Kappa yielded a model that almost never gave a perfect score. This led P1 to eventually adopt a performance metric better suited to their setting, which balanced predictive accuracy with some chance that students could get a perfect score. P1’s experience illustrates

that top-down, theoretical analysis of an established precedent (Quadratic Weighted Kappa) was insufficient for making a well-informed piggybacking decision. Instead, P1's experience illustrates the "science of the concrete" characteristic of bricolage, where tangible interactions with materials are important for understanding how to assemble elements in the repertoire into a functional artifact.

**4.2.2 Composing.** Several data scientists constructed prediction tasks by **composing**: *assembling multiple outcome variables into a single prediction target*. Often, data scientists used composing after realizing that a single outcome variable did not capture all dimensions of a construct they were trying to measure. For example, P8 used composing to predict the "academic work ethic" of college freshmen from their high school performance. While high school GPA was an established predictor of collegiate success, P8 felt that GPA alone would not offer enough specificity to differentiate work ethic from other success factors. P8 reasoned about how they might *compose* standardized test scores and GPA in a single measure for "academic work ethic" by combining the two:

*"If you assume that SAT or ACT scores are [...] indicators of academic knowledge, [...] this is where the person's at. If there's a difference between where the person should be and where they're at, we thought it was reasonable to conceptualize that as work ethic."*

P8 designed a regression model that predicted GPA given test scores, creating a model that accounted for the gap between them to capture "work ethic." A negative residual indicated that students had lower grades than predicted by their test scores, thereby signaling work ethic concerns. While teaching an introductory freshman class, P8 was surprised to find that students predicted to have poor work ethic were more willing than others to copy homework. Her model's predictions also closely aligned with students' long-term persistence outcomes (e.g., retention and graduation). These observations provided evidence for the *consequential validity* of their model's operationalization of "work ethic."

Yet some participants found that *composing* had simplicity limitations. For example, P2 tried to use a multi-outcome latent variable to predict student "engagement." However, when P2 presented the model to software engineers and product managers, they distrusted its predictions when the model contradicted their prior beliefs about new product features. P2 explained that:

*"If I show them a latent model, and it says, 'No, it doesn't look like you improved engagement' then there's gonna be a lot of questions about that model. What is this actually measuring? How do you explain it?"*

P2 found it easier to justify their model when it targeted a single outcome based on what software engineers and product managers viewed as "hard data." As a result, while P2 acknowledged validity benefits with latent variable models, they did not use them in industry practice. This tension between validity and simplicity highlights how bricolage solutions often bear the mark of compromise, as data scientists negotiate between theoretical ideals and practical constraints.

**4.2.3 Swapping.** As data scientists worked with their formulation, they sometimes reassessed their initial choice of outcome variable. The **swapping** strategy describes *changing the outcome variable when the initial one fails to meet evaluation criteria*. For example, P10 used swapping while developing a model to predict diabetes-related "health complications." P10 initially operationalized "health complications" using a cardiovascular outcome. Yet after discovering that cardiovascular outcomes were hard to predict, they later swapped to a measure of kidney function:

*"At some point, after [...] seeing that this [cardiovascular] condition is quite hard to predict, I thought that switching to something that is strongly related to the biomarkers I had available would make more sense. [...] Seeing the data, getting to work with data, you might need to change your outcomes."*

However, P10's use of renal outcomes was controversial because the renal biomarkers used as features were closely related to the outcome definition. As a result, it was not clear that the new model captured a true predictive signal in their data. This concern relates to *internal validity*, or "the existence of a defensible causal relationship between features and the target label" [28]. However, the goal of this project was to demonstrate a new prediction method, so P10 tolerated internal validity limitations. P10 explained that they would place greater weight on these internal validity concerns if they were deploying their model in a clinical context. P10's experience embodies the responsive dialogue central to bricolage, where "seeing the data, getting to work with data" led to a reconsideration of the initial approach based on the artifact that could be constructed with the materials at hand.

**4.2.4 Bridging.** While working with their data, two participants identified "gold standard" outcomes with established validity that were rarely available and infeasible to directly predict. Data scientists addressed this challenge by **bridging**: *targeting a low-cost, readily available outcome as a surrogate for the more expensive gold standard measure*. While bridging, participants evaluated their surrogate model using the gold standard outcome to construct a "bridge" between the surrogate model and the unobservable construct of interest.

For example, P6 used bridging to evaluate the "readability improvement" of content changes they made to an online learning platform. While P6 wanted to assess whether their changes helped struggling readers, they only had access to "gold standard" test scores from a small fraction of schools. P6 circumvented this issue by using students' performance on early lessons in the software as a surrogate for their reading ability:

*"We had this idea that came from another paper that [...] this early lesson [...] was a good proxy for reading ability. So we built some [...] models that took raw interaction data and [we] found a neural network model that does a reasonable job of predicting student's [test] performance from this first lesson."*

P6 validated their surrogate model by checking its performance against the test scores they *did* have available. Once they were confident that the model had "reasonable" validity, they used it to evaluate whether content changes helped students predicted to be poor readers. They found that the content updates offered a large readability improvement and deployed them on the platform, declaring the project "a really big success." P6 acknowledged that both the gold standard tests and the formulation were imperfect. However, given the analytics focus of the project, P6 had some tolerance for validity imperfections in their bridging formulation. This tolerance for imperfection, balanced against practical gains, exemplifies the bricoleur's pragmatic approach to problem-solving under resource constraints.

**4.2.5 Refining.** While constructing a target variable, participants also iteratively tweaked and refined their formulation by making small changes. This strategy involved **refining**: *modifying the target variable definition in response to a defect along one or more criteria*. In contrast to the larger change in outcomes involved with swapping, participants used refining to make more granular adjustments to their target variable.

For example, P13 used **refining** to narrow in on definitions for animal diseases in her dataset. While human healthcare had a mature ICD-9 disease classification scheme, P13 explained that animal health conditions remained vague and poorly specified. Clinicians were also unsure of how to map their understanding of pet diseases into a **label definition**: *or a systematic set of classification rules in data*. P13 addressed this issue by refining the label definition with clinicians.

*"Clinicians [...] can't think of how do I define the dog has diabetes? I would have the vet look at it and [ask if] there [are] symptoms and signs [...] I couldn't get right. That takes many times of iteration to find out."*

At each stage of refinement, P13 and the vet used written clinical notes about the pet to identify animals who were labeled incorrectly under the current definition. P13 then revised the label definition to correctly categorize the misclassified instances (e.g., adding a new rule-based white blood cell count). These iterations continued 20 to 30 times until the label definition sufficiently captured the disease. P13 emphasized a need for practicality during the process: they were fine with 10% of diseases being mislabeled if the major dimensions of the disease definition were covered. P13's acceptance of imperfection while ensuring "major dimensions" were covered typifies the bricoleur's focus on creating solutions that are "fit for purpose" rather than ideal.

**4.2.6 When Reformulation Fails.** After applying strategies, many data scientists eventually crafted a prediction task that satisfied their criteria. However, one discontinued a project after identifying a flaw that they could not patch through reformulation. P12 was developing a model to flag patients in the general care clinic who were at risk of Celiac disease. An initial model predicted test results with good accuracy. However, after carefully inspecting their formulation, P12 realized that patients with test results already appeared sick to physicians. Training a model on patients who *appeared sick* then deploying the model to *the full patient population* threatened the model's external validity. P12's team considered a variety of workarounds to patch this gap. However, when no solutions appeared feasible, P12 discontinued the project. While most data scientists applied strategies in response to visible defects in their formulation (e.g., predictive performance concerns), the flaw identified by P12 was *invisible*. In particular, identifying this issue required scrutinizing the connection between the dataset available for model development (i.e., patients who *appeared sick* to doctors) and the population at deployment time (i.e., the general population entering the clinic). We describe other approaches that data scientists used to shed light on invisible validity defects in Section 4.3.

### 4.3 Evaluating the Validity of a Problem Formulation

While applying strategies, data scientists often took a step back to evaluate the validity of their problem formulation. Drawing on domain knowledge, they engaged in theory-driven practices to probe conceptual relationships between their target variable and theoretical constructs related to their problem formulation. However, data scientists also treated their problem formulation as a material object to be scrutinized by inspecting concrete predictions generated on specific data instances. Lévi-Strauss calls this more tangible approach "the science of the concrete" — whereby the bricoleur understands their artifact through direct manipulation of material elements as opposed to abstract, conceptual reasoning. We now describe these theory-driven and data-driven practices in detail, as summarized in Table 3.

**4.3.1 Theory-Driven Evaluation Practices.** While engaging in theory-driven practices, data scientists used their domain knowledge to reason about a problem formulation *conceptually*, without scrutinizing a specific data sample. These theory-driven practices were informed by participants' prior modeling experience and an understanding of the sociotechnical context (e.g., hospitals, classrooms) in which their models operate. Though participants had varying levels of formal measurement theory training, many theory-driven practices connect to established construct validity sub-criteria.

While engaging with the problem formulation vignette, several participants identified spurious factors that might affect an outcome variable while remaining unrelated to the construct of interest. When considering "length of hospital stay" as a measure for "health acuity," P11 noted that "there's

*so many other things that could contribute to somebody being in the emergency department for 12 hours that might not have anything to do with how sick they are*”, citing staffing levels and emergency department crowding as examples. In validity theory, this reasoning connects to *divergent validity*, or the extent to which a problem formulation measures aspects of other unrelated constructs [28].

Several participants also used multiple outcomes to develop a more complete understanding of how a construct should be conceptualized. P2 wanted to understand whether students who used hints were “gaming” the lesson – i.e., using hints to finish questions without seriously engaging with the material. Accuracy did not provide enough information to identify gaming from hint use because students could answer a question incorrectly, even after using a hint and applying effort. However, P2 found that students took longer to respond when they were engaged with the lesson. As a result, a combination of *accuracy*, *response time*, and *hint use* was necessary to triangulate gaming behavior.<sup>4</sup> Reasoning across multiple outcomes connects to an assessment of content validity, or “the extent to which an operationalization wholly and fully captures the substantive nature of the construct purported to be measured” [62].

Finally, participants described their choice of outcome variable as closely linked to the intervention being informed by a predictive model. For example, P6 needed to decide which outcome time horizon (i.e., question lesson, week, semester, or year) was most appropriate for measuring the affect of platform content changes on student learning. P6 shared that “*to get in the right neighborhood, where would we reasonably expect to move the needle?*” If the intervention included many content changes, they leaned towards an end-of-year student performance outcome. However, if a change to the platform was more granular, they would start at a more fine grained question or lesson level. This assessment of the relationship between interventions and outcomes connects to *consequential validity*, or “the real-world consequences of using the measurements obtained from a measurement model” [62].

**4.3.2 Data-Driven Evaluation Practices.** While engaging in data-driven practices, participants checked whether their theoretical understanding of a construct was borne out by the measurements they observed on specific data instances. Gaps between theoretical expectations and observations signaled a potential flaw in their formulation to be scrutinized. By identifying, probing, and later reconciling these gaps through a strategy, data scientists slowly built trust in their formulation’s validity. In some cases, participants checked how a target variable categorized specific data instances (e.g., “ground truth” labels constructed by applying a label definition). In others, they checked the downstream predictions generated by a model trained on a target variable. Each of these practices depended heavily upon the specific data sample used to generate measurements. As such, our understanding of these activities emerged through the directed storytelling segment of interviews, as participants recounted specific times they manipulated a dataset in prior projects.

Several participants applied *spot checks and heuristics* to test whether measurements were inline with expected theoretical properties. For example, while evaluating new test items, P6 plotted students’ scores as a function of their number of question attempts. A “*nice, monotonically increasing curve*” with a “*logarithmic flavor*” provided evidence that the item encoded knowledge of a single skill. P9 checked that the base rate of a disease derived from diagnostic codes matched the expected population prevalence. A large mismatch was a signal that piggybacking off of existing disease definitions was infeasible. Participants described these heuristics as imperfect; a necessary but insufficient condition for demonstrating validity. These checks and heuristics are an example of

<sup>4</sup>Participants who reasoned across multiple outcomes sometimes suggested composing. However, in some cases, participants used multiple outcomes to evaluate a time-lagged prediction model without explicitly leveraging multiple outcomes in the predictive model when composing wasn’t feasible (e.g., due to simplicity concerns).

Type	Practice	Description	Example
Theory-Driven	<b>Identifying spurious causes of outcomes.</b>	Identifying factors influencing the outcome unrelated to the construct of interest.	<i>"There's so many other things that could contribute to somebody being in the emergency department for 12 hours that might not have anything to do with how sick they are." (P11)</i>
	<b>Reasoning across multiple outcomes.</b>	Leveraging multiple outcomes to gain a more holistic understanding of a construct.	<i>"If you have high response times, and they're using [learning] supports, that's a good indication that they're trying, even if they don't get the next problem correct." (P2)</i>
	<b>Linking interventions with outcomes.</b>	Reasoning about potential causal relationships between an intervention and outcomes.	<i>"To get in the neighborhood of the right outcome measure, it is often like, well, what does this intervention really look like? [...] Where would we ... expect to move the needle?" (P6)</i>
Data-Driven	<b>Validating predictions in real-world deployments.</b>	Checking the extent to which model predictions align with expectations in real-world conditions.	<i>"We did an evaluation study where the clinicians would actually review cases while the information that our models predicted would be relevant was highlighted for them." (P15)</i>
	<b>Poking holes in label definitions.</b>	Attempting to identify instances that are misclassified under a current label definition.	<i>"[It was] a lot of literally, in some cases, printing off student essays and handing them to employees that were former teachers and saying, why did this person give this label?" (P1)</i>
	<b>Applying quick checks and heuristics.</b>	Using pragmatic checks to quickly gauge whether the task setup matches the data sample.	<i>"There's several models that existed in the literature that I attempted to transfer to my work. I checked it [...] using finger in the wind sort of [approaches]. Like, what's the distribution of my data? Is it closer to 7% or is it closer to 20%?" (P7)</i>

Table 3. Data scientists engaged in both theory-driven and data-driven practices while evaluating validity.

evaluating *face validity*, or the degree to which an operationalization measures what it purports to measure based on a quick, superficial assessment [28].

Participants tried to *poke holes in label definitions* by identifying cases a label definition misclassified the instance in comparison to theoretical expectations. While *refining*, P13 tried to poke holes by identifying specific cases where the expert judgment of the clinician diverged from the categorization assigned by the label definition. In education, P1 tried to poke holes in essay scoring rubrics by printing off physical copies and asking employees who were former teachers to explain *"why did this person give this label?"* Employees' response served to build P1's intuitive sense of *how* the rubric was used to assign instance-level scores, and identify potential gaps. By *"staying close to the data"*, P1 had the trust needed in the labeling process to justify specific essay

scores to students, teachers, other external stakeholders. P1 described this trust as critical should a measurement decision later require formal legal justification — e.g., in the context of university admissions. For both P1 and P13, poking holes required working in tandem with domain experts, who had the knowledge required to verify whether the instance *really was* positive or negative for the construct of interest. The practice of poking holes and patching them connects to an assessment of *content validity*.

Finally, participants checked whether model behavior reflected their expectations in real-world deployment contexts. P15 checked the validity of a model trained to predict the “clinical relevance” of electronic medical records by conducting a field study that tested whether predictions saved clinicians time. Time savings provided evidence that the predicted fields really were “clinically relevant.” P15 and many other participants found both theory-driven and data-driven evaluation practices critical for establishing a holistic understanding of their formulation’s validity.

## 5 DISCUSSION

Today, significant focus across the HCI, CSCW, ML, and FAccT communities has turned towards evaluating the extent to which algorithmic systems achieve their purported aims. Given this attention, researchers and practitioners increasingly draw upon validity concepts established in the quantitative social sciences to describe desirable properties of algorithmic systems. For instance, a growing body of work examines how validity sub-criteria — such as content, convergent, and discriminant validity — might be applied to evaluate algorithmic systems [5, 28, 47, 62, 71, 83, 114, 115]. Yet our findings uncover a disconnect between the traditional, top-down measurement process assumed by validity theory and the ways data scientists navigate measurement decisions in practice. In particular, data scientists typically must make do with limited resources by re-appropriating existing data for new measurement tasks. These constraints force data scientists to balance validity with other important criteria, such as predictive performance, simplicity, and resource requirements.

Based on these findings, we argue that **the challenge is not to replace these bricolage practices with top-down planning, but rather to develop forms of scaffolding that can enhance the rigor of bricolage while preserving its creativity and adaptability**. Our findings provide a foundation for the research community to develop this scaffolding, by identifying concrete opportunities to help data scientists balance validity with other criteria, identify and apply problem (re)formulation strategies, and assess validity using both top-down and bottom-up evaluation practices. We now reflect on data science as a bricolage practice, before providing targeted recommendations for supporting data scientists’ problem formulation practices going forward.

### 5.1 Data Science as Bricolage

References to data science often invoke engineering metaphors: data scientists “engineer” features, “build” models, and “architect” computational solutions. Yet our findings suggest that target variable construction more closely resembles bricolage than engineering—a mode of knowledge production whereby the bricoleur makes do with the resources at hand rather than acquiring materials for a predefined plan. Our findings highlight how this process of making do unfolds in real-world projects, as data scientists re-appropriate data fields originally collected for entirely different purposes to accomplish their measurement goals. Bricolage describes *how* data scientists orient themselves towards measurement problems, as well as nature and quality of their solutions.

**5.1.1 The Data Science Repertoire.** The bricoleur operates with an inventory of odds and ends called a “repertoire.” Data scientists’ repertoire includes both material resources (e.g., datasets, computational tools) and immaterial resources (e.g., domain theories, methods) brought to bear

on a modeling task. The repertoire shapes the space of possibilities that bricoleurs imagine while confronting new problems. As Louridas [84] observes while drawing connections between bricolage and design, “the bricoleur asks his collection, whereas the engineer, like the scientist, asks the universe.” As one participant noted, a fundamental limitation is that data scientists cannot generate information that was never collected in the first place. The constrained nature of data scientists’ repertoire forces them to re-examine their limited materials from new perspectives.

In particular, materials in data scientists’ repertoire function as what Lévi-Strauss calls “signs”, or concrete objects that represent abstract concepts, but remain constrained by their imagined uses. P6’s creative re-purposing of raw user interactions in an online learning platform as a proxy for “reading ability” exemplifies the use of a sign, by unexpectedly transforming a mundane data field into an operationalization of a construct. Yet, as Levi-Strauss [81] notes, “the possibilities remain always limited by the particular history of each piece, and by what is predetermined in it due to the original usage for which it was conceived.” Had P6 *not* been aware of this proxy for “reading ability”, established through a chance encounter with a colleague’s work, they may have never made this connection.

**5.1.2 Engaging in Dialogue with Data.** After taking stock of the repertoire, the bricoleur engages in dialogue with their materials to identify a workable solution to the problem at hand. As Lévi-Strauss explains, the bricoleur must “engage into a kind of dialog with [the repertoire], to index, before choosing among them, the possible answers that the set can offer to his problem.” [81] This dialogue not only shapes the bricoleur’s understanding of each element in isolation, but also their interactions: “one element’s possibilities interact with all other elements’ possibilities, with the overall organization of the artifact he makes. The results of these interactions are never what he expects, and he must respond to them.” [84] P10’s experience illustrates how dialogue unfolds in data science. After initially selecting a cardiovascular outcome based on precedent in academic literature (piggybacking), “*seeing the data, getting to work with data*” led them to *swap* to a renal outcome that was more predictable given available biomarkers. This example illustrates problem (re)formulation strategies serve as the brushstroke of measurement in data science work. Like a painter with a canvas, data scientists apply a strategy, step back to survey their progress along multiple criteria, then consider a new strategy in response to newly-identified flaws.

Dialogue marks a notable shift from the top-down planning approach to measurement. Under top-down planning, a practice better described as a *monologue* [135], the data scientist defines a theoretical construct of interest, designs an operationalization, *then* collects data. Yet this monologue supposes that the data scientist can navigate the space of constraints and contingencies in their formulation through mental simulation, without a specific data sample or model. Our findings underscore that this is rarely the case. Instead, data scientists required *tangible interactions* with data to understand its affordances and constraints. It was through this very process of tangible interaction with data, borne out of resource constraints, that data scientists made creative leaps in their projects. As P6 explained, had a “magical oracle” granted their team access to all the data they required, they would have simply collected a comprehensive dataset with student test outcomes. It was the combination of resource constraints and tangible interaction with data that drove P6’s team to adaptively *bridge* using the limited data at hand.

**5.1.3 The Emerging Artifact.** Critically, the artifact resulting from bricolage is rarely ideal, but rather bears the marks of the constraints, compromises, and creator that shaped it. As Louridas [84] explains, “bricolage is … at the mercy of contingencies, either external, in the form of influences, constraints, and adversities of the external world, or internal, in the form of the creator’s idiosyncrasy.” Analogously, data scientists readily acknowledged imperfection in their solutions. P13 explained that a label definition that miscategorized 10% of pet diseases was acceptable if the

“major dimensions” of the disease’s theoretical conceptualization were covered. P6 deemed their re-purposing of existing data a “really big success”, but also identified important limitations. Thus, from the bricoleur’s perspective, the key question is not whether a solution is *ideal*, but rather whether it is *fit for purpose*. P6 embodied this ethos when they justified *bridging* by explaining that their platform analytics use case was “not high stakes.” Had their model been designed for use cases they perceived as higher-stakes, P6 explained that they may have applied a higher level of stringency.

Beyond constraints, Lévi-Strauss notes a personal dimension to the artifact that emerges from bricolage: “he [the bricoleur] ‘speaks’ [...] through the choices he makes among the limited possibilities, the character and the life of the creator.” P8’s creative use of *composing*, which conceptualizes “academic work ethic” as the difference between *where a student is at* and *where they should be*, may have never emerged had they not also observed real students in real classrooms as an educator. Similarly, P12’s project, which had strong buy-in from multiple stakeholders and momentum towards completion, never came to fruition due to a validity flaw spotted by P12. Yet it was through P12’s personal interest in adjacent bodies of biostatistics literature, and their attention to methodological rigor in problem formulation, that they became aware of this otherwise invisible flaw. Each of these examples illustrates that the journey through forking paths in problem formulation depends greatly on the knowledge and experience of the bricoleur.

In sum, bricolage acknowledges that data scientists must work within the confines of their existing resources and organizational processes. These practical considerations are neither centered nor prioritized in the conventional top-down approach to measurement. Thus, short of dramatically expanding the available resources for a project, a rigid enforcement of top-down planning is unlikely to translate to a desirable project outcome. **Therefore, rather than recommending more stringent enforcement of the traditional top-down measurement approach, we suggest more effectively scaffolding existing bricolage practices.** We now revisit traditional top-down measurement interventions with this insight in mind, before envisioning new ways of scaffolding data science practice.

## 5.2 Design Implications: Reconciling Top-Down Measurement Interventions with Bricolage Practice

Our findings uncover a tension between existing interventions designed to help analysts implement the traditional top-down measurement workflow, and data scientists’ bottom-up bricolage practices. However, this tension is not unique to data science. History reveals similar patterns across domains where top-down planning approaches have failed while bottom-up alternatives succeeded. In urban planning, modernist urban renewal projects that imposed regimented public works projects onto complex social systems often destroyed functioning neighborhoods, while community-based planning approaches showed strong success [40, 124]. In software development, rigid and systematic waterfall planning methodologies gave way to agile approaches that embrace iteration and adaptation [103]. Key to the success of these bottom-up alternatives was not wholesale rejection of structure, but rather targeted interventions like mixed-use zoning in urban planning and sprint retrospectives in agile development. These interventions recognized that while existing practices were in need of reform, effective solutions must work within resource constraints and pinpoint specific procedural weaknesses in need of constructive redirection. Similarly, when improving measurement practices in data science, interventions tailored to support data scientists’ *existing* practices are likely to be more effective than direct translation of rigid, top-down measurement frameworks. In the following sections, we examine how existing top-down measurement interventions might be adapted to better support data scientists’ bricolage practices.

**5.2.1 Pre-Registration for Predictive Modeling.** Pre-registration is one top-down intervention that is widely recommended in the quantitative social sciences [138]. This protocol discourages data-dependent decision-making by requiring an analyst to specify outcome measures and hypotheses *before* running analyses. Recent work argues for the adoption of pre-registration in predictive modeling contexts [58] – i.e., by requiring data scientists to finalize their problem formulation before training a model. Yet requiring such pre-registration is akin to *blindfolding the bricoleur* – it restricts the very bottom-up, iterative process that data scientists require to make do under resource constraints. Indeed, data scientists interviewed by Hofman et al. [58] worried that pre-registration might limit their creativity, or make it challenging to later tweak their formulation if they started off with limited familiarity with their data.

A more targeted protocol for limiting data-driven decision-making might encourage teams to establish minimum performance standards for each criterion before beginning problem formulation. This approach could help to prevent *standard erosion* – or a selective reduction in stringency applied to a criterion over subsequent rounds of (re)formulation. Our interviews surfaced evidence for standard erosion in several projects, as participants weakened validity standards to improve predictive performance or simplicity. To mitigate standard erosion, HCI and CSCW researchers might co-design tools such as radar plot visualizations (e.g., Figure 3) to help teams chart trade-off spaces across criteria and identify suitable operating regions. For instance, a research team might prioritize validity over simplicity, while an educational technology company might be more willing to accept validity trade-offs for improved simplicity. By making performance trade-offs explicit, teams can (1) prevent standard erosion, (2) enable iterative (re)formulation, and (3) provide rationale for discontinuing projects when minimum standards cannot be satisfied with available resources.

**5.2.2 Latent Variable Models.** Researchers in the quantitative social sciences have also developed *statistical tools* to improve the validity of models. For example, latent variable models (e.g., [10, 14]) use statistical machinery to characterize the relationship between unobservable theoretical constructs and observed proxies. Recent work has proposed latent variable models (e.g., [91]) and related modeling advances (e.g., [46, 140]) to improve the validity of models used for decision support and recommendation systems. Yet these sophisticated approaches require advanced statistical training, multiple outcome variables, and significant implementation time. This imposes practical barriers that impede adoption. For example, P1 shared that, while they were aware of advanced multi-outcome modeling approaches, the “extra bells and whistles” provided by these approaches made them less reliable in large-scale deployment contexts.

This disconnect between theoretical benefits of latent variable models and practical constraints presents an opportunity to develop modeling tools that better align with bricolage practices. Rather than introducing increasingly complex modeling tools, future work might focus on *simplifying* latent variable models to preserve their validity benefits while reducing implementation complexity. These simplified approaches could emphasize interpretability and robustness, helping teams effectively communicate their models to stakeholders without advanced statistical training. We elaborate further on improved modeling approaches designed to better facilitate bricolage in Section 5.3.3.

**5.2.3 Pedagogical Interventions.** Our findings also highlight opportunities for pedagogical interventions to better equip practitioners to navigate target variable construction. A key insight is that data scientists often draw on prior experience, engaging in analogical reasoning by identifying parallels between their current project and those encountered previously [64]. This highlights the value of **case-based learning** in data science education. Given the improvisational nature of real-world data science practice, case-based approaches enable students to analyze concrete case studies and understand how experienced data scientists make trade-offs between criteria like validity, simplicity, and resource requirements. Such case-based approaches could build upon

existing protocols designed to help practitioners identify potential failure modes of AI systems during the early stages of AI system design [66, 78, 121].

**Reflective practice** and **hands-on workshops** offer complementary roles in developing data scientists' capacity to reason about measurement and validity. This mirrors how tasks like feature engineering are often described—as a “craft”—involving domain expertise, theoretical grounding, and direct engagement with data [35]. Workshops focused on (re)formulation strategies—such as piggybacking, composing, swapping, bridging, and refining—can expose students to the tactical, trial-and-error process of shaping a prediction task within real-world constraints. Together, these approaches offer an integrated method for teaching the complexities of bricolage-based problem formulation.

### 5.3 Design Implications: Developing a Scaffolding for Effective Bricolage Practice

In addition to re-visiting existing top-down interventions, our findings surface a need for entirely *new* interventions designed to help data scientists engage in bricolage more effectively. While many toolkits exist to help data scientists engage in inner-loop activities, such as training models [21, 33, 105, 146] and exploring data [74, 144], few support specific bricolage practices highlighted in our findings. Data scientists currently lack tools to help them weigh validity with competing criteria in their formulation, identify and apply strategies, and evaluate validity sub-criteria. **We now use our findings as a concrete scaffolding for identifying actionable points of intervention in data scientists' bricolage practices.** These interventions are intended to *facilitate* specific practices that data scientists found critical during problem formulation, while *adding guardrails around others* that have the potential to introduce flaws.

**5.3.1 Helping Data Scientists Weigh Tradeoffs.** In the quantitative social sciences, validity is traditionally viewed as the primary criterion for the scientific integrity of a study. However, while engaging in bricolage, data scientists were forced to adaptively balance the validity of their formulation with competing criteria. Yet this balancing act remains poorly supported, leaving the (re)formulation processes susceptible to inadvertent introduction of validity defects.

First, our findings illustrate a need to help data scientists design measures for their formulation's simplicity, validity, portability, and resource requirements. Teams currently rely heavily on measurements of predictive performance (e.g., AU-ROC or F1 scores) while applying strategies. Echoing the classic adage “what gets measured matters” [130], this practice leads teams to prioritize predictive performance improvements at the expense of other criteria. To address this imbalance, HCI and CSCW researchers might work with teams to co-design rubrics that help them operationalize measures for specific validity sub-criteria (e.g., convergent, divergent, content), simplicity (e.g., model complexity, implementation requirements), and portability (e.g., data field availability across contexts). For instance, a team might use such a rubric to set acceptable inter-rater reliability thresholds between automated essay scoring and human expert ratings. To support ease of implementation, this rubric could be paired with existing Responsible AI toolkits designed to help teams weigh the suitability of predictive modeling tools for specific applications [28, 71].

Second, there is an opportunity to develop tools that help data scientists identify targeted data collection opportunities. While engaging in bricolage, data scientists were forced to weigh the predictive performance and validity benefits of data collection against its resource requirements. Yet existing model training and evaluation toolkits often assume that datasets are fixed [12, 21, 146]. Future HCI and CSCW work might close this gap by developing tools that help data scientists quantify the improvements offered by collecting additional data. The framework of data valuation offers tools for ascribing monetary value to data resources [23, 109]. Future toolkits could build upon these frameworks to enable teams to estimate the benefits of new data sources prior to acquisition,

potentially leveraging synthetic data generation to quickly assess potential improvements before collecting real data. These efforts would be enhanced by empirical studies investigating how teams acquire and ascribe worth to data throughout the (re)formulation process.

**5.3.2 Helping Data Scientists Identify Strategies.** Our findings present an opportunity to help teams *identify* strategies during the target variable construction process. Data scientists currently identify strategies in an ad hoc, unstructured manner – e.g., by drawing upon their bespoke knowledge of academic literature, or the expertise of others on their team. This approach led teams to miss opportunities to improve their formulation. For instance, following the conclusion of their project, P12 identified a modeling framework—i.e., case-control methods [122]—that was tailor-made to address their project’s data constraints. Yet P12 was not aware of this framework because their background in machine learning did not involve training in biostatistics methods. This presents an opportunity to broaden data scientists’ awareness of relevant strategies before helping them narrow in on which are best suited to their bespoke modeling challenges.

Therefore, future HCI and CSCW research might develop systems that help teams exchange knowledge surrounding common (re)formulation challenges and mitigation strategies. While community-based knowledge exchange is central to inner-loop workflows, as illustrated by Stack Overflow [11] for coding and Hugging Face [128] for model training, data scientists’ support for problem formulation remains siloed. Similar to other online communities, this system might structure discussions via a comment forum where teams share specific strategies—P8’s *combining* strategy for measuring “academic work ethic” or P6’s *bridging* strategy for measuring “reading ability.” This forum could eventually be used to inform the design of consensus-backed protocols for navigating common (re)formulation challenges.

**5.3.3 Helping Data Scientists Apply Strategies.** Beyond strategy identification, our findings also highlight opportunities to help data scientists apply each strategy more effectively. Therefore, we now describe opportunities to develop tools and processes that support data scientists in applying each strategy while avoiding potential limitations.

**Piggybacking.** While *piggybacking*, data scientists port an existing problem formulation to their own modeling context. Piggybacking saves time when there is a strong alignment between established precedent and data scientists’ modeling context. Yet, as illustrated by P1’s experiences, piggybacking can also introduce serious defects when contextual differences are overlooked. Going forward, there is a need for structured protocols that help teams thoughtfully reason about the transferability of problem formulations between modeling contexts. While many guidelines support dataset and model documentation (e.g., [28, 42, 92]), none support assessments of problem formulation transferability. HCI and CSCW researchers might fill this gap by working with data scientists to co-design a rubric that teams can complete while considering piggybacking. This rubric might guide data scientists through a side-by-side comparison of an existing formulation and their own along dimensions such as high-level modeling goals, available data, and rationale for a selected target variable. Helping data scientists spot transferability gaps can also help teams provide a rationale to external stakeholders should a team decide to use a formulation that deviates from an established precedent.

**Composing.** While *composing*, data scientists combine multiple outcomes into a single target variable. Composing can improve validity by helping data scientists capture all relevant dimensions of a construct in a target variable. Yet composing can also introduce simplicity limitations. This tension indicates a need for simple and interpretable mechanisms for combining outcomes. Future HCI and ML research might explore opportunities for interpretable multi-outcome models that cleanly map conceptual relationships (e.g., “*the difference between where the person should be and where they’re at*”, P8) into transparent model designs. One approach likely to be effective

involves training multiple models to target multiple outcomes of interest (e.g., cost of medical care, re-admission), then aggregating predictions into a single score at run-time. P11 described such an approach as common in the medical industry, where both simplicity and validity are paramount. Several participants also suggested this approach during the design probe. HCI, ML, and Visualization researchers can implement this approach by developing tools that help teams encode their domain knowledge into aggregation functions that weight multiple predictions into a single score.

**Swapping.** While *swapping*, data scientists switch to a different outcome after they identify defects in one or more criteria. Swapping can be adaptive if it helps teams identify a formulation that satisfies all minimum performance standards. Yet swapping can also introduce standard erosion. This tension highlights a need for tools to help data scientists make more systematic swapping decisions. While many existing systems support feature engineering (e.g., [13, 48, 68, 116]) and EDA (e.g., [74, 144]), only one supports *outcome exploration* [41]. HCI and Visualization researchers can support outcome exploration by developing toolkits that help data scientists weigh the portability, simplicity, predictive performance, and resource requirements of alternative outcomes. Because our results suggest that swapping can be enticing given its predictive performance benefits, it is especially important that toolkits help teams identify validity drawbacks of swapping decisions. We discuss opportunities to support validity evaluation in Section 5.3.4.

**Bridging.** While *bridging*, data scientists adopt a low-cost, readily available target variable as a proxy for a more costly gold standard measure. When used effectively, bridging is a strong example of bricolage because it involves a creative approach for making do with limited data. Yet the success of bridging hinges on data scientists' ability to identify trustworthy sources of auxiliary information. To support bridging, data scientists need tools to help them identify, evaluate, and incorporate relevant sources of auxiliary data. While a large body of HCI and Data Engineering research has developed tools supporting data discovery (e.g., [15, 39, 101, 108]), existing toolkits are not tailored to support bridging. Future research might close this gap by developing tools that help data scientists forage for relevant auxiliary data and evaluate its expected utility against the time and cost required to incorporate it in their formulation. As with swapping, it is important that such tools help data scientists identify any defects associated with bridging.

**Refining.** While *refining*, data scientists make granular adjustments to how a target variable is coded in their data. Participants often engaged in refining by modifying a label definition used to operationalize a construct using a set of indicators used to categorize instances. For instance, P13 engaged in refining by working with domain experts to spot gaps in label definitions, then patch them with a revised set of indicators. Though refining can improve validity, practitioners lack tools to encode domain expert's knowledge into label definitions. Instead, existing tools developed by the HCI and ML communities help annotators apply a *fixed* label definition to annotate data (e.g., [30, 32, 96, 100]) or use crowdsourcing to refine annotation guidelines (e.g., [17, 24, 25, 67, 86]).<sup>5</sup> Future HCI and ML research can help data scientists encode domain experts' knowledge in label definitions by developing tools that help teams collaboratively systematize concepts [3] by iteratively identifying and refining indicators used in label definitions. Because our results indicate that domain experts require rich contextual information to identify areas where the systemization of a concept is incomplete, tools supporting this process should capture rich contextual information (e.g., written notes) in data.

### 5.3.4 Helping Data Scientists Evaluate Validity.

Disciplines in the quantitative social sciences have established numerous guidelines for evaluating the validity of a measurement instrument.

---

<sup>5</sup>Traditional annotation guidelines are *de facto* label definitions if they serve as a measurement instrument that operationalizes a latent construct of interest.

Such guidelines are often grouped into validity sub-criteria — such as convergent, divergent, and predictive validity — describing abstract, theoretical assertions about the expected behavior of a measurement instrument [34]. Yet our findings illustrate a disconnect between these abstract assertions and the more concrete theory- and data-driven evaluation activities performed by practitioners. We resolve this tension by identifying opportunities to scaffold data scientists' existing evaluation practices (Table 3).

First, our findings highlight opportunities to help teams incorporate domain knowledge while engaging in theory-driven evaluation activities. For example, HCI, ML, and Visualization researchers might develop toolkits that help teams map hypothesized causal relationships between predictive features, interventions, outcome variables, and latent constructs during problem formulation. In an education context, this tool might help teams chart students' longitudinal trajectories through the education system, weigh points of intervention, and consider an intervention's effect on multiple downstream outcomes of interest. This tool might also encourage teams to draw upon their domain knowledge to identify spurious causes of outcome variables — e.g., by listing causes of an outcome that are unrelated to the construct of interest. During the design probe, several participants independently suggested an interest in a tool supporting these theory-driven evaluation activities. While the *DoWhy* package [127] enables data scientists to reason about causal relationships in their data, this package is not tailored to support the formulation of prediction tasks.

Our findings also demonstrate a need for tools to help teams translate their theoretical understanding of a construct to low-level assertions about expected model behavior. *Behavioral evaluation*, which involves testing the capabilities of a system against a specification of requirements [22, 113], is one paradigm for facilitating this translation. Prior HCI research has explored leveraging behavioral evaluation to help data scientists evaluate the accuracy and fairness of models [21]. Future research might extend this approach to facilitate improved measurement practices. For instance, such a toolkit might enable data scientists to verify the scores assigned by a predictive model against *anchor points*, or cases known to be positive or negative for the construct of interest. When anchor points are unavailable, teams might also check consistency of predictions against a partial ordering over cases — e.g., student essays sorted by a domain expert for their expected “authenticity.” Supporting both theory- and data-driven activities is critical for helping teams close the gap between their abstract understanding of a construct and its operationalization in a modeling formulation.

#### 5.4 Implications for Generative AI Evaluation

While we study measurement in the context of predictive modeling, our findings also connect to emerging discourse surrounding the evaluation of Generative AI systems. Recent work has cast Generative AI evaluation as a measurement task, in which evaluation designers measure properties of AI systems — such as the “toxicity” or “helpfulness” of outputs — by (1) constructing a systematic definition for a concept to be measured, (2) operationalizing this definition into a measurement instrument (e.g., annotation instructions), and (3) applying this operationalization to obtain measurements [139]. Recent empirical studies have also shown that the evaluation design process is deeply iterative: practitioners often require multiple rounds of refinement to reconcile their conceptual understanding of evaluation criteria with measurements obtained on specific model outputs [8, 55, 104, 126, 133]. Teams may balance a range of practical criteria while engaging with this process, such as the extensibility and actionability of measurement instruments [53]. More broadly, viewing Generative AI evaluation design as a bricolage practice opens a space of research questions. How do evaluation designers make do when balancing validity with competing criteria under resource constraints (e.g., limited human ratings and real-world usage data)? How might

we develop a scaffolding for this iterative process that enhances validity, while also supporting flexibility? Our work offers an entry point to these key questions.

### 5.5 Study Limitations

Our study has several important limitations. First, by foregrounding data scientists' perspectives, we capture only one dimension of a deeply collaborative process. During interviews, data scientists described collaborating with domain experts, such as clinicians or educators, to refine label definitions and evaluate the validity of their formulations. Recent work has also shown that community members provide critical feedback on problem formulations [50, 78]. As a result, the factors shaping how bricolage unfolds likely extend beyond participants' vantage points as data scientists.

Our study also characterizes how problem formulation unfolds in a limited range of domains. We interviewed participants working in education and healthcare — where specific regulatory constraints, data ecosystems, and incentives may have shaped their bricolage practices. Data scientists in domains with more readily available data (e.g., social media platforms), different regulatory environments (e.g., finance), or distinct organizational incentives (e.g., user engagement versus educational assessment) may consider criteria beyond those documented in our findings.

Finally, our findings draw upon data scientists with a baseline level of domain knowledge and technical competency. Our recruitment criteria stipulated that data scientists have at least one year of full-time experience in the healthcare or education sectors to be eligible to participate. Additionally, our quoted examples of (re)formulation strategies in Section 4.2 often feature creative practices by highly experienced practitioners. Less experienced data scientists may draw upon a narrower range of strategies in their repertoire or engage in less sophisticated validity evaluation practices. Understanding and supporting novice practitioners' bricolage practices is an important area for future research.

## 6 CONCLUSION

In this study, we interviewed fifteen data scientists working in education and healthcare domains to understand their practices, challenges, and perceived opportunities for target variable construction in predictive modeling. We explore how data scientists adopt a *bricolage* approach to target variable construction, designing operationalizations for unobserved, latent constructs by *making do* with the data that they *already* have at their disposal. We characterize the contours and constraints of data scientists' bricolage practices, including the problem (re)formulation strategies they apply to balance across different criteria. Critically, we argue that *interventions* designed to improve target variable construction should begin by acknowledging the inherent uncertainty and resource constraints involved in real-world data science work. By understanding and designing to support target variable construction as a *bricolage process*—for example, by helping teams navigate validity tradeoffs, determine when and how to collect more data, and identify appropriate (re)formulation strategies—the research community can help data scientists more thoughtfully engage with problem formulation and avert the negative consequences that arise from poor modeling decisions.

## 7 ACKNOWLEDGMENTS

We thank our research participants for making this work possible and anonymous reviewers for their thoughtful feedback. This work was supported by an award from the UL Research Institutes through the Center for Advancing Safety of Machine Intelligence (CASMI) at Northwestern University and the National Science Foundation Graduate Research Fellowship Program (Award No. DGE1745016).

## REFERENCES

- [1] 2016. <https://www2.ed.gov/rschstat/eval/high-school/early-warning-systems-brief.pdf>
- [2] Amina A Abdu, Irene V Pasquetto, and Abigail Z Jacobs. 2023. An empirical analysis of racial categories in the algorithmic fairness literature. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 1324–1333.
- [3] Robert Adcock and David Collier. 2001. Measurement validity: A shared standard for qualitative and quantitative research. *American political science review* 95, 3, 529–546.
- [4] Yongsu Ahn and Yu-Ru Lin. 2019. Fairsight: Visual analytics for fairness in decision making. *IEEE transactions on visualization and computer graphics* 26, 1, 1086–1095.
- [5] Ahmed Alaa, Thomas Hartvigsen, Niloufar Golchini, Shiladitya Dutta, Frances Dean, Inioluwa Deborah Raji, and Travis Zack. 2025. Medical Large Language Model Benchmarks Should Prioritize Construct Validity. *arXiv preprint arXiv:2503.10694* (2025).
- [6] Sara Alspaugh, Nava Zokaei, Andrea Liu, Cindy Jin, and Marti A Hearst. 2018. Futzling and moseying: Interviews with professional data analysts on exploration practices. *IEEE transactions on visualization and computer graphics* 25, 1, 22–31.
- [7] Audrey Amrein-Beardsley. 2014. Rethinking value-added models in education: Critical perspectives on tests and assessment-based accountability. Routledge.
- [8] Zahra Ashktorab, Michael Desmond, Qian Pan, James M Johnson, Martin Santillan Cooper, Elizabeth M Daly, Rahul Nair, Tejaswini Pedapati, Swapna Achintalwar, and Werner Geyer. 2024. Aligning Human and LLM Judgments: Insights from EvalAssist on Task-Specific Evaluations and AI-assisted Assessment Strategy Preferences. *arXiv preprint arXiv:2410.00873* (2024).
- [9] Ted Baker and Reed E Nelson. 2005. Creating something from nothing: Resource construction through entrepreneurial bricolage. *Administrative science quarterly* 50, 3, 329–366.
- [10] David J Bartholomew, Martin Knott, and Irini Moustaki. 2011. *Latent variable models and factor analysis: A unified approach*. John Wiley & Sons.
- [11] Anton Barua, Stephen W Thomas, and Ahmed E Hassan. 2014. What are developers talking about? an analysis of topics and trends in stack overflow. *Empirical software engineering* 19, 619–654.
- [12] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilović, et al. 2019. AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development* 63, 4/5, 4–1.
- [13] Aditya Bhattacharya, Simone Stumpf, Lucija Gosak, Gregor Stiglic, and Katrien Verbert. 2024. EXMOS: Explanatory Model Steering Through Multifaceted Explanations and Data Configurations. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–27.
- [14] Christopher M Bishop. 1998. Latent variable models. In *Learning in graphical models*. Springer, 371–403.
- [15] Alex Bogatu, Alvaro AA Fernandes, Norman W Paton, and Nikolaos Konstantinou. 2020. Dataset discovery in data lakes. In *2020 ieee 36th international conference on data engineering (icde)*. IEEE, 709–720.
- [16] Geoffrey C Bowker. 2000. *Sorting Things Out: Classification and Its Consequences*. MIT press.
- [17] Jonathan Bragg, Mausam, and Daniel S Weld. 2018. Sprout: Crowd-powered task design for crowdsourcing. In *Proceedings of the 31st annual acm symposium on user interface software and technology*. 165–176.
- [18] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative research in psychology* 3, 2, 77–101.
- [19] Monika Büscher, Satinder Gill, Preben Mogensen, and Dan Shapiro. 2001. Landscapes of practice: Bricolage as a method for situated design. *Computer Supported Cooperative Work (CSCW)* 10, 1–28.
- [20] Ángel Alexander Cabrera, Will Epperson, Fred Hohman, Minsuk Kahng, Jamie Morgenstern, and Duen Horng Chau. 2019. FairVis: Visual analytics for discovering intersectional bias in machine learning. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*. IEEE, 46–56.
- [21] Ángel Alexander Cabrera, Erica Fu, Donald Bertucci, Kenneth Holstein, Ameet Talwalkar, Jason I Hong, and Adam Perer. 2023. Zeno: An interactive framework for behavioral evaluation of machine learning. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–14.
- [22] Ángel Alexander Cabrera, Marco Tulio Ribeiro, Bongshin Lee, Robert Deline, Adam Perer, and Steven M Drucker. 2023. What did my AI learn? How data scientists make sense of model behavior. *ACM Transactions on Computer-Human Interaction* 30, 1, 1–27.
- [23] Raul Castro Fernandez. 2023. Data-sharing markets: model, protocol, and algorithms to incentivize the formation of data-sharing consortia. *Proceedings of the ACM on Management of Data* 1, 2, 1–25.
- [24] Joseph Chee Chang, Saleema Amershi, and Ece Kamar. 2017. Revolt: Collaborative crowdsourcing for labeling machine learning datasets. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 2334–2346.

- [25] Quan Ze Chen and Amy X Zhang. 2023. Judgment sieve: Reducing uncertainty in group judgments through interventions targeting ambiguity versus disagreement. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW2, 1–26.
- [26] Hao-Fei Cheng, Logan Stapleton, Anna Kawakami, Venkatesh Sivaraman, Yanghuidi Cheng, Diana Qing, Adam Perer, Kenneth Holstein, Zhiwei Steven Wu, and Haiyi Zhu. 2022. How child welfare workers reduce racial disparities in algorithmic decisions. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–22.
- [27] Victoria Clarke and Virginia Braun. 2017. Thematic analysis. *The journal of positive psychology* 12, 3, 297–298.
- [28] Amanda Coston, Anna Kawakami, Haiyi Zhu, Ken Holstein, and Hoda Heidari. 2023. A validity perspective on evaluating the justified use of data-driven decision-making algorithms. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*. IEEE, 690–704.
- [29] Alfred W Crosby. 1997. *The measure of reality: Quantification in Western Europe, 1250-1600*. Cambridge University Press.
- [30] Tobias Daudert. 2020. A web-based collaborative annotation and consolidation tool. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. 7053–7059.
- [31] Wesley Hanwen Deng, Manish Nagireddy, Michelle Seng Ah Lee, Jatinder Singh, Zhiwei Steven Wu, Kenneth Holstein, and Haiyi Zhu. 2022. Exploring how machine learning practitioners (try to) use fairness toolkits. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 473–484.
- [32] Michael Desmond, Michael Muller, Zahra Ashktorab, Casey Dugan, Evelyn Duesterwald, Kristina Brimijoin, Catherine Finegan-Dollak, Michelle Brachman, Aabhas Sharma, Narendra Nath Joshi, et al. 2021. Increasing the speed and accuracy of data labeling through an ai assisted interface. In *Proceedings of the 26th International Conference on Intelligent User Interfaces*. 392–401.
- [33] Dennis Dingen, Marcel van't Veer, Patrick Houthuizen, Eveline HJ Mestrom, Erik HHM Korsten, Arthur RA Bouwman, and Jarke Van Wijk. 2018. RegressionExplorer: Interactive exploration of logistic regression models with subgroup analysis. *IEEE transactions on visualization and computer graphics* 25, 1, 246–255.
- [34] Ellen A Drost. 2011. Validity and reliability in social science research. *Education Research and perspectives* 38, 1, 105–123.
- [35] Pablo Duboue. 2020. *The art of feature engineering: essentials for machine learning*. Cambridge University Press.
- [36] Raffi Duymedjian and Charles-Clemens Rüling. 2010. Towards a foundation of bricolage in organization and management theory. *Organization studies* 31, 2, 133–151.
- [37] Shelley Evenson. 2016. Driving Service Design By Directed Storytelling. In *Design for services*. Routledge, 66–72.
- [38] B Everett. 2013. *An introduction to latent variable models*. Springer Science & Business Media.
- [39] Raul Castro Fernandez, Ziawasch Abedjan, Famien Koko, Gina Yuan, Samuel Madden, and Michael Stonebraker. 2018. Aurum: A data discovery system. In *2018 IEEE 34th International Conference on Data Engineering (ICDE)*. IEEE, 1001–1012.
- [40] Martin Fuller and Ryan Moore. 2017. *An Analysis of Jane Jacobs's The Death and Life of Great American Cities*. Macat Library.
- [41] Dalia Gala, Milo Phillips-Brown, Naman Goel, Carinal Prunkl, Laura Alvarez Jubete, Ray Eitel-Porter, et al. 2024. FairTargetSim: An Interactive Simulator for Understanding and Explaining the Fairness Effects of Target Variable Definition. *arXiv preprint arXiv:2403.06031*.
- [42] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. Datasheets for datasets. *Commun. ACM* 64, 12, 86–92.
- [43] Darren Gergle and Desney S Tan. 2014. Experimental research in HCI. In *Ways of Knowing in HCI*. Springer, 191–227.
- [44] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. 2000. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *circulation* 101, 23, e215–e220.
- [45] Nitesh Goyal, Ian D Kivlichan, Rachel Rosen, and Lucy Vasserman. 2022. Is your toxicity my toxicity? exploring the impact of rater identity on toxicity annotation. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2, 1–28.
- [46] Luke Guerdan, Amanda Coston, Kenneth Holstein, and Zhiwei Steven Wu. 2023. Counterfactual prediction under outcome measurement error. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 1584–1598.
- [47] Luke Guerdan, Amanda Coston, Zhiwei Steven Wu, and Kenneth Holstein. 2023. Ground (less) truth: A causal framework for proxy labels in human-algorithm decision-making. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 688–704.
- [48] Philip J Guo, Sean Kandel, Joseph M Hellerstein, and Jeffrey Heer. 2011. Proactive wrangling: Mixed-initiative end-user programming of data transformation scripts. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*. 65–74.

- [49] Ian Hacking. 1999. *The social construction of what?* Harvard university press.
- [50] MD Romael Haque, Devansh Saxena, Katy Weathington, Joseph Chudzik, and Shion Guha. 2024. Are we asking the right questions?: Designing for community stakeholders' interactions with ai in policing. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–20.
- [51] Bill Harding. 2021. *Software effort estimates vs popular developer productivity metrics*. Technical Report. GitClear.
- [52] Emma Harvey, Hauke Sandhaus, Abigail Z Jacobs, Emanuel Moss, and Mona Sloane. 2024. The Cadaver in the Machine: The Social Practices of Measurement and Validation in Motion Capture Technology. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–23.
- [53] Emma Harvey, Emily Sheng, Su Lin Blodgett, Alexandra Chouldechova, Jean Garcia-Gathright, Alexandra Olteanu, and Hanna Wallach. 2025. Understanding and Meeting Practitioner Needs When Measuring Representational Harms Caused by LLM-Based Systems. *arXiv preprint arXiv:2506.04482* (2025).
- [54] Orit Hazzan and Jim Tomayko. 2004. Human aspects of software engineering: The case of extreme programming. In *International Conference on Extreme Programming and Agile Processes in software Engineering*. Springer, 303–311.
- [55] Zeyu He, Saniya Naphade, and Ting-Hao Kenneth Huang. 2025. Prompting in the Dark: Assessing Human Performance in Prompt Engineering for Data Labeling When Gold Labels Are Absent. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–33.
- [56] Douglas D Heckathorn. 2011. Comment: Snowball versus respondent-driven sampling. *Sociological methodology* 41, 1, 355–366.
- [57] Nicole Hoess, Carlos Paradis, Rick Kazman, and Wolfgang Mauerer. 2025. Does the Tool Matter? Exploring Some Causes of Threats to Validity in Mining Software Repositories. *arXiv preprint arXiv:2501.15114* (2025).
- [58] Jake M Hofman, Angelos Chatzimpampas, Amit Sharma, Duncan J Watts, and Jessica Hullman. 2023. Pre-registration for predictive modeling. *arXiv preprint arXiv:2311.18807*.
- [59] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudik, and Hanna Wallach. 2019. Improving fairness in machine learning systems: What do industry practitioners need?. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–16.
- [60] Youyang Hou and Dakuo Wang. 2017. Hacking with NPOs: collaborative analytics and broker roles in civic data hackathons. *Proceedings of the ACM on Human-Computer Interaction* 1, CSCW, 1–16.
- [61] Steven J Ingels, Daniel J Pratt, James E Rogers, Peter H Siegel, and Ellen S Stutts. 2004. Education Longitudinal Study of 2002: Base Year Data File User's Manual. NCES 2004-405. *National Center for Education Statistics*.
- [62] Abigail Z Jacobs and Hanna Wallach. 2021. Measurement and fairness. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 375–385.
- [63] Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. 2020. Mimic-iv. *PhysioNet*. Available online at: <https://physionet.org/content/mimiciv/1.0/> (accessed August 23, 2021), 49–55.
- [64] Eunice Jun, Melissa Birchfield, Nicole De Moura, Jeffrey Heer, and Rene Just. 2022. Hypothesis formalization: Empirical findings, software limitations, and design implications. *ACM Transactions on Computer-Human Interaction (TOCHI)* 29, 1 (2022), 1–28.
- [65] Eunice Jun, Audrey Seo, Jeffrey Heer, and René Just. 2022. Tisane: Authoring statistical models via formal reasoning from conceptual and data relationships. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [66] Ji-Young Jung, Devansh Saxena, Minjung Park, Jini Kim, Jodi Forlizzi, Kenneth Holstein, and John Zimmerman. 2025. Making the Right Thing: Bridging HCI and Responsible AI in Early-Stage AI Concept Selection. *arXiv preprint arXiv:2506.17494* (2025).
- [67] V K. Chaithanya Manam, Dwarakanath Jampani, Mariam Zaim, Meng-Han Wu, and Alexander J. Quinn. 2019. Taskmate: A mechanism to improve the quality of instructions in crowdsourcing. In *Companion Proceedings of The 2019 World Wide Web Conference*. 1121–1130.
- [68] Sean Kandel, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. 2011. Wrangler: Interactive visual specification of data transformation scripts. In *Proceedings of the sigchi conference on human factors in computing systems*. 3363–3372.
- [69] Sean Kandel, Andreas Paepcke, Joseph M Hellerstein, and Jeffrey Heer. 2012. Enterprise data analysis and visualization: An interview study. *IEEE transactions on visualization and computer graphics* 18, 12, 2917–2926.
- [70] Sayash Kapoor and Arvind Narayanan. 2023. Leakage and the reproducibility crisis in machine-learning-based science. *Patterns* 4, 9 (2023).
- [71] Anna Kawakami, Amanda Coston, Haiyi Zhu, Hoda Heidari, and Kenneth Holstein. 2024. The Situate AI Guidebook: Co-Designing a Toolkit to Support Multi-Stakeholder, Early-stage Deliberations Around Public Sector AI Proposals. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–22.
- [72] Anna Kawakami, Venkatesh Sivaraman, Hao-Fei Cheng, Logan Stapleton, Yanghuidi Cheng, Diana Qing, Adam Perer, Zhiwei Steven Wu, Haiyi Zhu, and Kenneth Holstein. 2022. Improving human-AI partnerships in child welfare: understanding worker practices, challenges, and desires for algorithmic decision support. In *Proceedings of the 2022*

- CHI Conference on Human Factors in Computing Systems.* 1–18.
- [73] Mary Beth Kery and Brad A Myers. 2017. Exploring exploratory programming. In *2017 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. IEEE, 25–29.
  - [74] Mary Beth Kery, Marissa Radensky, Mahima Arya, Bonnie E John, and Brad A Myers. 2018. The story in the notebook: Exploratory data science using a literate programming tool. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–11.
  - [75] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2016. The emerging role of data scientists on software development teams. In *Proceedings of the 38th International Conference on Software Engineering*. 96–107.
  - [76] Amy J Ko, Robin Abraham, Laura Beckwith, Alan Blackwell, Margaret Burnett, Martin Erwig, Chris Scaffidi, Joseph Lawrence, Henry Lieberman, Brad Myers, et al. 2011. The state of the art in end-user software engineering. *ACM Computing Surveys (CSUR)* 43, 3, 1–44.
  - [77] Sean Kross and Philip Guo. 2021. Orienting, framing, bridging, magic, and counseling: How data scientists navigate the outer loop of client collaborations in industry and academia. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2, 1–28.
  - [78] Tzu-Sheng Kuo, Hong Shen, Jisoo Geum, Nev Jones, Jason I Hong, Haiyi Zhu, and Kenneth Holstein. 2023. Understanding Frontline Workers' and Unhoused Individuals' Perspectives on AI Used in Homeless Services. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17.
  - [79] Bruno Latour and Steve Woolgar. 2013. *Laboratory life: The construction of scientific facts*. Princeton university press.
  - [80] Jason Lefever, Yuanfang Cai, Humberto Cervantes, Rick Kazman, and Hongzhou Fang. 2021. On the lack of consensus among technical debt detection tools. In *2021 IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 121–130.
  - [81] Claude Levi-Strauss. 1966. *The savage mind*. University of Chicago Press.
  - [82] Jessica Liu, Huaming Chen, Jun Shen, and Kim-Kwang Raymond Choo. 2024. FairCompass: Operationalising fairness in machine learning. *IEEE Transactions on Artificial Intelligence*.
  - [83] Yu Lu Liu, Su Lin Blodgett, Jackie Chi Kit Cheung, Q Vera Liao, Alexandra Olteanu, and Ziang Xiao. 2024. ECBD: Evidence-centered benchmark design for NLP. *arXiv preprint arXiv:2406.08723* (2024).
  - [84] Panagiotis Louridas. 1999. Design as bricolage: Anthropology meets design thinking. *Design Studies* 20, 6 (1999), 517–535.
  - [85] Scott Lundberg. 2017. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*.
  - [86] VK Manam and Alexander Quinn. 2018. Wingit: Efficient refinement of unclear task instructions. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 6. 108–116.
  - [87] Yaoli Mao, Dakuo Wang, Michael Muller, Kush R Varshney, Ioana Baldini, Casey Dugan, and Aleksandra Mojsilović. 2019. How data scientists work together with domain experts in scientific collaborations: To find the right answer or to ask the right question? *Proceedings of the ACM on Human-Computer Interaction* 3, GROUP, 1–23.
  - [88] Sara Mateus and Soumodip Sarkar. 2024. Bricolage—a systematic review, conceptualization, and research agenda. *Entrepreneurship & Regional Development*, 1–22.
  - [89] Jennifer Mickel. 2024. Racial/Ethnic Categories in AI and Algorithmic Fairness: Why They Matter and What They Represent. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 2484–2494.
  - [90] Delbert C Miller and Neil J Salkind. 2002. *Handbook of research design and social measurement*. Sage.
  - [91] Smitha Milli, Luca Belli, and Moritz Hardt. 2021. From optimizing engagement to measuring value. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 714–722.
  - [92] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*. 220–229.
  - [93] A Mol. 2002. *The body multiple: Ontology in medical practice*. Duke University Press.
  - [94] Sendhil Mullainathan and Ziad Obermeyer. 2021. On the inequity of predicting A while hoping for B. In *AEA Papers and Proceedings*, Vol. 111. American Economic Association 2014 Broadway, Suite 305, Nashville, TN 37203, 37–42.
  - [95] Michael Muller, Ingrid Lange, Dakuo Wang, David Piorkowski, Jason Tsay, Q Vera Liao, Casey Dugan, and Thomas Erickson. 2019. How data science workers work with data: Discovery, capture, curation, design, creation. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–15.
  - [96] Michael Muller, Christine T Wolf, Josh Andres, Michael Desmond, Narendra Nath Joshi, Zahra Ashktorab, Aabhas Sharma, Kristina Brimijoin, Qian Pan, Evelyn Duesterwald, et al. 2021. Designing ground truth and the social life of labels. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–16.
  - [97] Alexander Martin Mussgnug. 2022. The predictive reframing of machine learning applications: good predictions and bad measurements. *European Journal for Philosophy of Science* 12, 3 (2022), 55.
  - [98] Brad A Myers, Amy J Ko, Thomas D LaToza, and YoungSeok Yoon. 2016. Programmers are users too: Human-centered methods for improving programming tools. *Computer* 49, 7, 44–52.

- [99] Nadia Nahar, Shurui Zhou, Grace Lewis, and Christian Kästner. 2022. Collaboration challenges in building ml-enabled systems: Communication, documentation, engineering, and process. In *Proceedings of the 44th international conference on software engineering*. 413–425.
- [100] Hiroki Nakayama, Takahiro Kubo, Junya Kamura, Yasufumi Taniguchi, and Xu Liang. 2018. doccano: Text Annotation Tool for Human. <https://github.com/doccano/doccano> Software available from <https://github.com/doccano/doccano>.
- [101] Fatemeh Nargesian, Erkang Zhu, Renée J Miller, Ken Q Pu, and Patricia C Arocena. 2019. Data lake management: challenges and opportunities. *Proceedings of the VLDB Endowment* 12, 12, 1986–1989.
- [102] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. 2019. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366, 6464, 447–453.
- [103] Jorge A Osorio, Michel RV Chaudron, and Werner Heijstek. 2011. Moving from waterfall to iterative development: An empirical evaluation of advantages, disadvantages and risks of RUP. In *2011 37th EUROMICRO Conference on Software Engineering and Advanced Applications*. IEEE, 453–460.
- [104] Qian Pan, Zahra Ashktorab, Michael Desmond, Martin Santillan Cooper, James Johnson, Rahul Nair, Elizabeth Daly, and Werner Geyer. 2024. Human-Centered Design Recommendations for LLM-as-a-judge. *arXiv preprint arXiv:2407.03479* (2024).
- [105] Bo Pang, Erik Nijkamp, and Ying Nian Wu. 2020. Deep learning with tensorflow: A review. *Journal of Educational and Behavioral Statistics* 45, 2, 227–248.
- [106] Samir Passi and Solon Barocas. 2019. Problem formulation and fairness. In *Proceedings of the conference on fairness, accountability, and transparency*. 39–48.
- [107] Samir Passi and Steven J Jackson. 2018. Trust in data science: Collaboration, translation, and accountability in corporate data science projects. *Proceedings of the ACM on human-computer interaction* 2, CSCW, 1–28.
- [108] Norman W Paton, Jiaoyan Chen, and Zhenyu Wu. 2023. Dataset discovery and exploration: A survey. *Comput. Surveys* 56, 4, 1–37.
- [109] Juan Carlos Perdomo. 2023. The Relative Value of Prediction in Algorithmic Decision Making. *arXiv preprint arXiv:2312.08511*.
- [110] Kathleen H Pine and Max Liboiron. 2015. The politics of measurement and action. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. 3147–3156.
- [111] Theodore M Porter. 1996. *Trust in numbers: The pursuit of objectivity in science and public life*. Princeton University Press.
- [112] Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. 2022. Data cards: Purposeful and transparent dataset documentation for responsible ai. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1776–1826.
- [113] Iyad Rahwan, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia Breazeal, Jacob W Crandall, Nicholas A Christakis, Iain D Couzin, Matthew O Jackson, et al. 2019. Machine behaviour. *Nature* 568, 7753, 477–486.
- [114] Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. 2021. AI and the everything in the whole wide world benchmark. *arXiv preprint arXiv:2111.15366* (2021).
- [115] Inioluwa Deborah Raji, I Elizabeth Kumar, Aaron Horowitz, and Andrew Selbst. 2022. The fallacy of AI functionality. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 959–972.
- [116] Tye Rattenbury, Joseph M Hellerstein, Jeffrey Heer, Sean Kandel, and Connor Carreras. 2017. *Principles of data wrangling: Practical techniques for data preparation*. " O'Reilly Media, Inc".
- [117] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1135–1144.
- [118] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [119] Pedro Saleiro, Benedict Kuester, Loren Hinkson, Jesse London, Abby Stevens, Ari Anisfeld, Kit T Rodolfa, and Rayid Ghani. 2018. Aequitas: A bias and fairness audit toolkit. *arXiv preprint arXiv:1811.05577*.
- [120] Benjamin Saunders, Julius Sim, Tom Kingstone, Shula Baker, Jackie Waterfield, Bernadette Bartlam, Heather Burroughs, and Clare Jinks. 2018. Saturation in qualitative research: exploring its conceptualization and operationalization. *Quality & quantity* 52, 1893–1907.
- [121] Devansh Saxena, Ji-Youn Jung, Jodi Forlizzi, Kenneth Holstein, and John Zimmerman. 2025. AI Mismatches: Identifying Potential Algorithmic Harms Before AI Development. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–23.
- [122] James J Schlesselman. 1982. *Case-control studies: design, conduct, analysis*. Vol. 2. Oxford university press.
- [123] Moritz Schubotz, Ankit Satpute, André Greiner-Petter, Akiko Aizawa, and Bela Gipp. 2022. Caching and Reproducibility: Making Data Science Experiments Faster and FAIRer. *Frontiers in Research Metrics and Analytics* 7 (2022),

861944.

- [124] James C Scott. 2020. *Seeing like a state: How certain schemes to improve the human condition have failed*. yale university Press.
- [125] Ahmed Seffah, Jan Gulliksen, and Michel C Desmarais. 2005. *Human-centered software engineering-integrating usability in the software development lifecycle*. Vol. 8. Springer Science & Business Media.
- [126] Shreya Shankar, JD Zamfirescu-Pereira, Björn Hartmann, Aditya Parameswaran, and Ian Arawjo. 2024. Who validates the validators? aligning llm-assisted evaluation of llm outputs with human preferences. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*. 1–14.
- [127] Amit Sharma and Emre Kiciman. 2020. DoWhy: An end-to-end library for causal inference. *arXiv preprint arXiv:2011.04216*.
- [128] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yuetong Zhuang. 2024. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems* 36.
- [129] Kultar Singh. 2007. *Quantitative social research methods*. Sage.
- [130] Wesley G Skogan. 1999. Measuring what matters. In *Proceedings from the policing research institute meetings*. Department of Justice, 37–88.
- [131] Jonathan Stray, Ivan Vendrov, Jeremy Nixon, Steven Adler, and Dylan Hadfield-Menell. 2021. What are you optimizing for? aligning recommender systems with human values. *arXiv preprint arXiv:2107.10939*.
- [132] Lucy Suchman. 1993. Do categories have politics? The language/action perspective reconsidered. *Computer supported cooperative work (CSCW)* 2, 177–190.
- [133] Annalisa Szymanski, Simret Araya Gebreegziabher, Oghenemaro Anuyah, Ronald A Metoyer, and Toby Jia-Jun Li. 2024. Comparing Criteria Development Across Domain Experts, Lay Users, and Models in Large Language Model Evaluation. *arXiv preprint arXiv:2410.02054* (2024).
- [134] Florian Tramer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, Jean-Pierre Hubaux, Mathias Humbert, Ari Juels, and Huang Lin. 2017. Fairtest: Discovering unwarranted associations in data-driven applications. In *2017 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 401–416.
- [135] Sherry Turkle. 2011. *Life on the Screen*. Simon and Schuster.
- [136] Bill Turque. 2012. Creative... motivating'and fired. *The Washington Post* 6.
- [137] Anna Vallgårda and Ylva Fernaeus. 2015. Interaction design as a bricolage practice. In *Proceedings of the ninth international conference on tangible, embedded, and embodied interaction*. 173–180.
- [138] Anna Elisabeth Van't Veer and Roger Giner-Sorolla. 2016. Pre-registration in social psychology—A discussion and suggested template. *Journal of experimental social psychology* 67, 2–12.
- [139] Hanna Wallach, Meera Desai, Nicholas Pangakis, A Feder Cooper, Angelina Wang, Solon Barocas, Alexandra Chouldechova, Chad Atalla, Su Lin Blodgett, Emily Corvi, et al. 2024. Evaluating Generative AI Systems is a Social Science Measurement Challenge. *arXiv preprint arXiv:2411.10939* (2024).
- [140] Jamelle Watson-Daniels, Solon Barocas, Jake M Hofman, and Alexandra Chouldechova. 2023. Multi-Target Multiplicity: Flexibility and Fairness in Target Specification under Resource Constraints. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 297–311.
- [141] Hilde Weerts, Miroslav Dudík, Richard Edgar, Adrin Jalali, Roman Lutz, and Michael Madaio. 2023. Fairlearn: Assessing and improving fairness of ai systems. *Journal of Machine Learning Research* 24, 257, 1–8.
- [142] James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viégas, and Jimbo Wilson. 2019. The what-if tool: Interactive probing of machine learning models. *IEEE transactions on visualization and computer graphics* 26, 1, 56–65.
- [143] Kanit Wongsuphasawat, Yang Liu, and Jeffrey Heer. 2019. Goals, process, and challenges of exploratory data analysis: An interview study. *arXiv preprint arXiv:1911.00568*.
- [144] Kanit Wongsuphasawat, Dominik Moritz, Anushka Anand, Jock Mackinlay, Bill Howe, and Jeffrey Heer. 2015. Voyager: Exploratory analysis via faceted browsing of visualization recommendations. *IEEE transactions on visualization and computer graphics* 22, 1, 649–658.
- [145] Feng Xie, Jun Zhou, Jin Wee Lee, Mingrui Tan, Siqi Li, Logasan S/O Rajnthern, Marcel Lucas Chee, Bibhas Chakraborty, An-Kwok Ian Wong, Alon Dagan, et al. 2022. Benchmarking emergency department prediction models with machine learning and public electronic health records. *Scientific Data* 9, 1, 658.
- [146] Jing Nathan Yan, Ziwei Gu, and Jeffrey M Rzeszotarski. 2021. Tessera: Discretizing data analysis workflows on a task level. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–15.
- [147] Amy X Zhang, Michael Muller, and Dakuo Wang. 2020. How do data science workers collaborate? roles, workflows, and tools. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW1, 1–23.

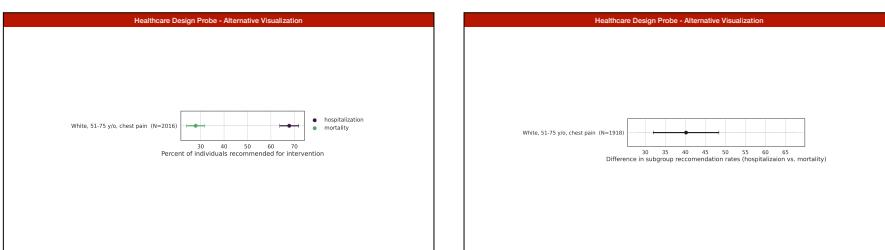
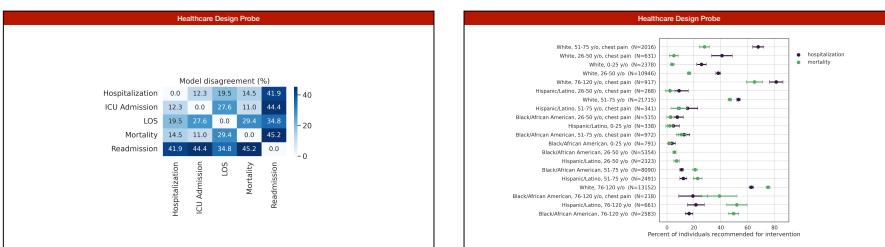
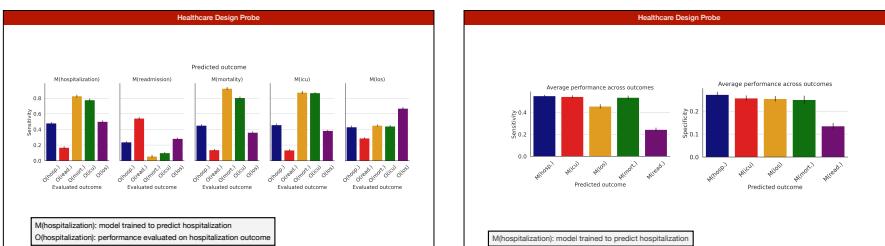
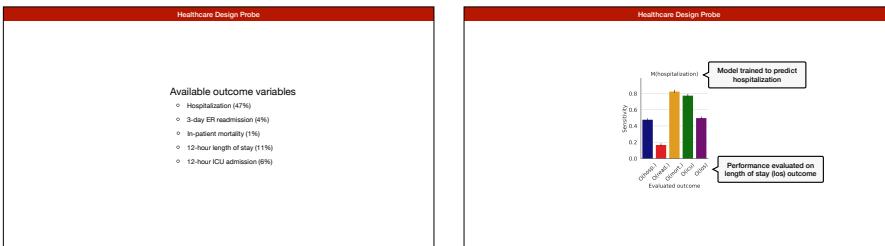
## A DETAILS OF THE VIGNETTE-BASED PROBLEM FORMULATION TASK

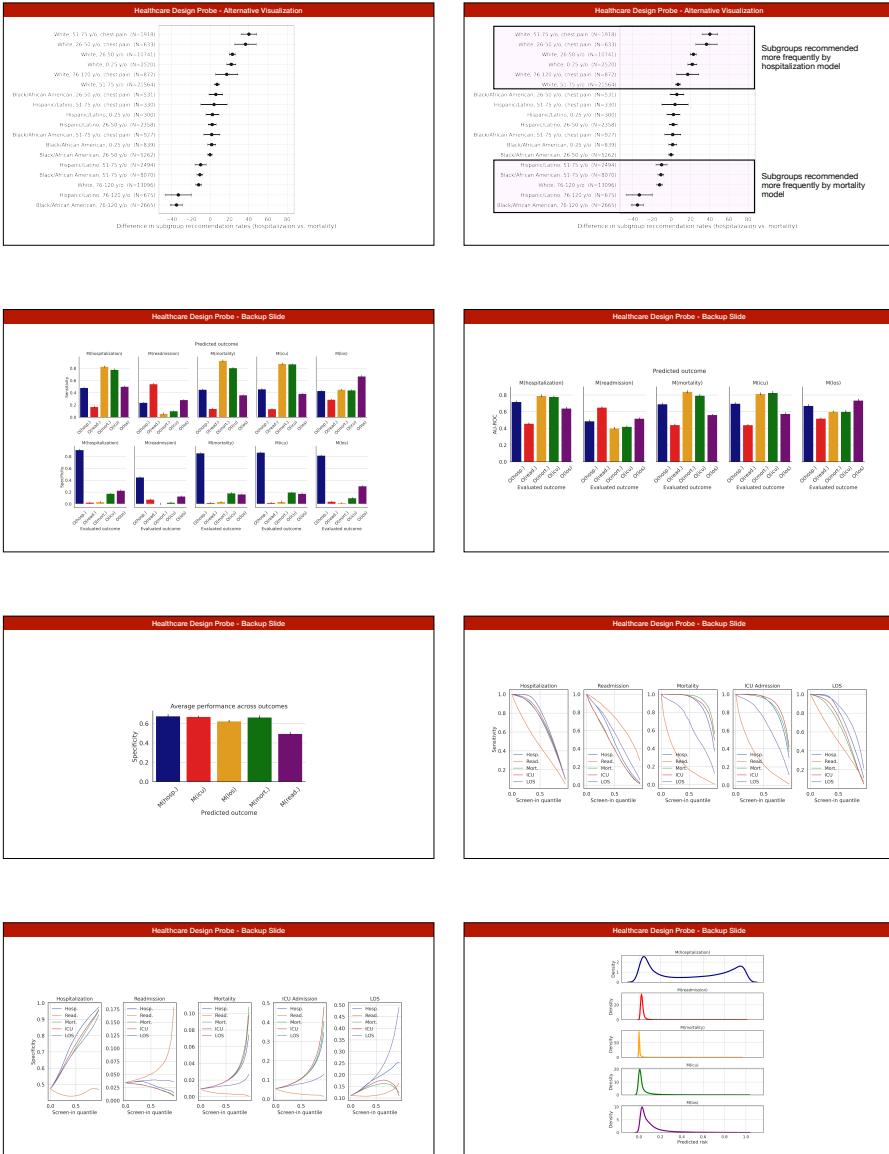
We began the design activity by introducing participants to a modeling scenario matched to their area of domain expertise. We asked participants working in healthcare to imagine that they were developing a model to support emergency room triage decisions by predicting patients with high “*medical acuity*.” We told participants that patients’ electronic medical records and current symptoms were available as predictors, and described five downstream health outcomes that were available as prediction targets: hospitalization, 3-day ER readmission, in-patient mortality, 12-hour length of stay, and 12-hour ICU admission. To understand how participants assessed the *face validity* of alternative formulations, we provided them with the base rate of each outcome and asked them to reason about which might serve as “*better or worse*” measures of acuity given the stated modeling goals. Data for this vignette was drawn from MIMIC-IV ED [44, 63], a dataset commonly used to benchmark predictive models targeting the same set of outcome variables (e.g., hospitalization, in-patient mortality, etc.) in the machine learning for healthcare literature [145].

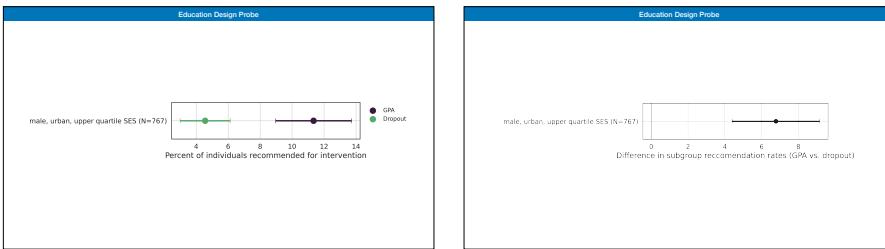
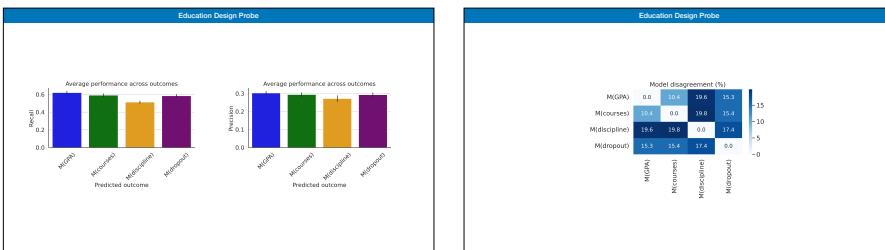
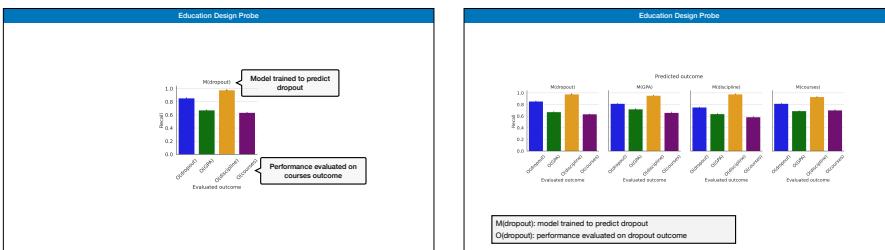
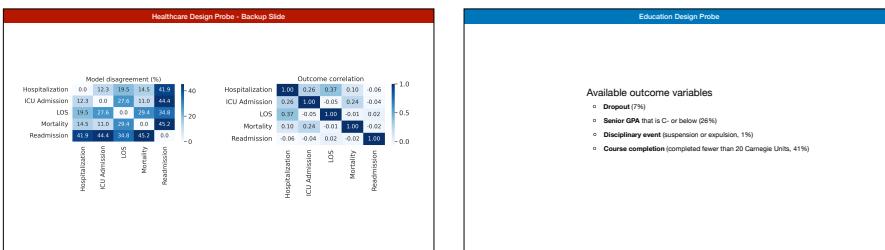
We asked participants working in education to imagine that they were developing a model to be used as part of an Early Warning System (EWS) for high school sophomores. We told participants that this system would be used to match students predicted to be “at-risk” with a personalized academic success intervention. Following US Department of Education guidelines [1], we defined “at-risk” students as those failing to achieve basic proficiency in core academic subjects (e.g., reading, math). We told participants that they had access to predictors such as students’ academic history, extracurricular involvement, and demographic factors, in addition to four outcomes recorded at the end of their senior year: dropout, senior GPA at a C- or below, major disciplinary event (i.e., suspension or expulsion), and course completion. To understand how participants weighted *face validity*, we provided them with the base rate of each outcome and asked them to reason about which might serve as “*better or worse*” measures of academic risk given the stated modeling goals. We leveraged data from the Educational Longitudinal Study (ELS) [61] to construct the vignette.

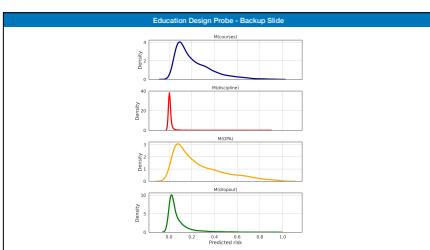
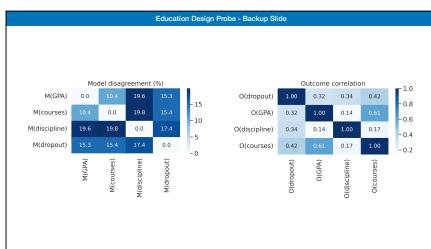
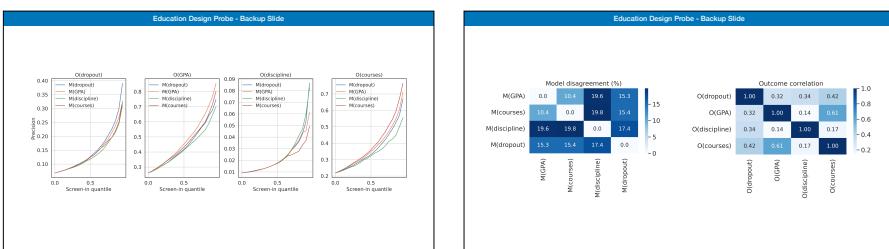
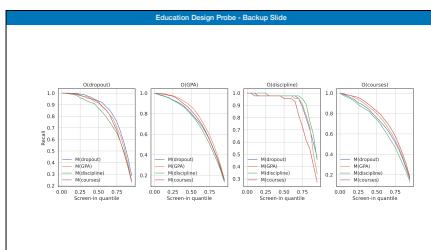
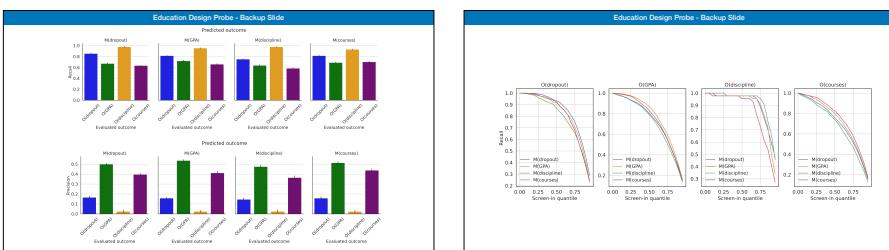
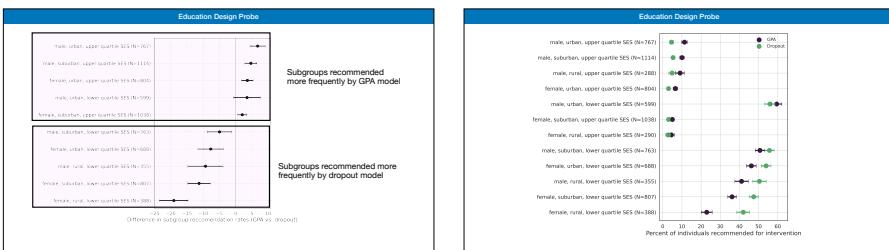
After introducing the vignette, we presented a series of model evaluations indicating the performance of predictive models targeting each outcome. After showing each plot, we paused and asked participants to reflect on which outcomes might serve as “*better or worse choices*” for the stated modeling goals. We also asked participants to share any additional information that would be useful for informing their thought process. We provide screenshots for each evaluation plot shown to participants in the pages included below.

We piloted our protocol on one participant to check its flow and comprehensibility. Then, as non-pilot participants engaged with the task over subsequent interviews, we refined the specificity of the scenario description and accompanying evaluation plots to build upon participants’ questions and feedback. For example, we later added reports of models’ sensitivity and specificity, AU-ROC curves, and kernel density plots after several participants reported interest in this information. Furthermore, while the initial vignette offered general descriptions of modeling intervention (i.e., instructing them to “*decide which patients should receive additional medical resources*” or “*identify high school sophomores in need of student success interventions*”), we later increased the specificity of the intervention description (provided in bullet points above) to enable participants to reason about the formulation in greater granularity.









Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009