

The Necessity and Difficulty of Navigating Uncertainty to Develop an Individual-Level Computational Model*

Alexander J. Freund¹[0000-0002-9791-9245] and Philippe J. Giabbanelli¹[0000-0001-6816-355X]

Department of Computer Science & Software Engineering, Miami University, Oxford OH 45056, USA <https://www.dachb.com> {freundaj,giabbapj}@miamioh.edu

Abstract. The design of an individual-level computational model requires modelers to deal with uncertainty by making assumptions on causal mechanisms (when they are insufficiently characterized in a problem domain) or feature values (when available data does not cover all features that need to be initialized in the model). The simplifications and judgments that modelers make to construct a model are not commonly reported or rely on evasive justifications such as ‘for the sake of simplicity’, which adds another layer of uncertainty. In this paper, we present the first framework to transparently and systematically investigate which factors should be included in a model, where assumptions will be needed, and what level of uncertainty will be produced. We demonstrate that it is computationally prohibitive (i.e. NP-Hard) to create a model that supports a set of interventions while minimizing uncertainty. Since heuristics are necessary, we formally specify and evaluate two common strategies that emphasize different aspects of a model, such as building the ‘simplest’ model in number of rules or actively avoiding uncertainty.

Keywords: Agent-Based Model · Causal Map · Graph Algorithms · Information Fusion.

1 Introduction

The design of an individual-level model (e.g., Agent-Based Model) is a common activity in computational science. In computational *social* science, such models can serve to explain social phenomena or safely test interventions within an artificial society before selecting the most promising ones for real-world pilot testing [13, 14, 7]. In an individual-level model, the simulated entities have their own *features* (e.g., age, gender, beliefs and values) which need to be initialized at the start of the simulation (i.e. given a *baseline value*). As the simulation unfolds, the entities’ behaviors and some features will be updated based on a set of *rules*

* Alexander J. Freund thanks the Department of Computer Science & Software Engineering and the Graduate School at Miami University for research funding. Both authors thank Ketra Rice and Nisha Nataraj for providing a case study.

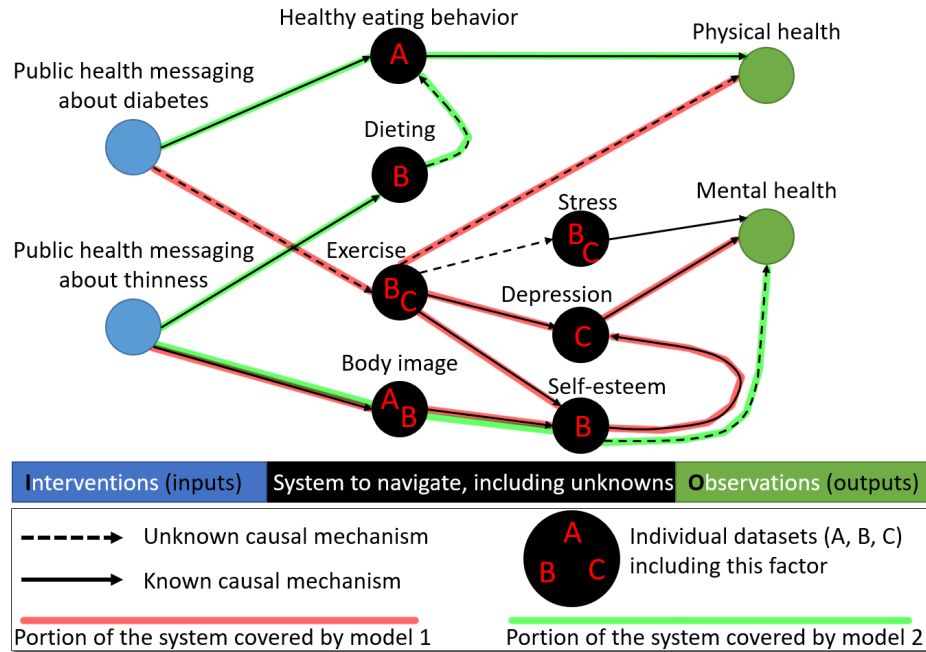


Fig. 1. A systems map (nodes and directed edges) of a problem space can be used to identify what salient factors and rules will go into a model based on available data and a systematic handling of uncertainties.

that can take into account the entity’s features as well as the features of simulated peers or the environment. By specifying the rules, a modeler expresses how to build reactive entities that either continue to engage in an existing pattern of behavior or adopt a new one by reacting to socio-environmental stimuli. The recommended guidance is that a “model should be embedded in existing theories and make use of whatever data are available. [Its] assumptions need to be clearly articulated, supported by the existing theories and justified by whatever information is available.” [1]

In contrast to the next steps which are more standardized (e.g., implementation in an object-oriented language, verification and validation) [11, 3, 10], the design of a model and the extent to which it should use information or theories is particularly subject to ad-hoc practices, hence to variations across modeling teams. To illustrate these variations, consider Figure 1 which exemplifies the problem space (i.e., set of relevant factors and interrelationships) in the case of food-related behaviors. For a model to be ‘fit for purpose’ (i.e. adequate) [34], the rules should connect the interventions required by model users to observable model outcomes. In this situation, two teams could produce and justify very different models. One team with access to dataset *B* could create model 1 (Fig 1, red highlights), which focuses on exercise and also touches on depression, body

image, and self-esteem. This model will require making 3 assumptions: one to compensate for the lack of baseline data to initialize the entities' depression level, and two for unknown rules relating public messaging to exercise, or exercise to physical health. Another team with access to dataset *A* could create model 2 (Fig 1, green highlights) which is more focused on eating and also requires three assumptions. A team with access to *both* datasets *A* and *B* could further reduce the number of assumptions necessary, as they could pass on representing dieting or the link from self-esteem to mental health while keeping the model fit for purpose. The design of an individual-based model is thus heavily impacted by how a team handles *parameter uncertainty* either in unknown causal mechanisms or in unknown feature values. However, the “different types of simplifications and scientific judgments that have to be made” [2] are not commonly reported, in part since they are not required by standardized documentations [16] (e.g., ODD, ODD+D). As essential aspects of model building are shaped by important yet unknown decisions, we face the additional problem of *structural uncertainty* [2].

Although the need to lower structural uncertainty in model building has long been established and mentioned in methodological guidance [35], this task has been hampered by the lack of a clear picture on the problem space (i.e. the map shown in Figure 1). In other words, it's simple to *state* that ‘modelers should explain how they systematically navigate the uncertainties in the problem space based on available data and assumptions’, but it's difficult for modelers to follow a *systematic method* to handle a problem space that has not been precisely mapped. Indeed, the creation of a comprehensive map of a problem space¹ (known as a ‘systems map’) is often beyond the scope of creating one model, which usually only involves a brief literature review [1] and/or consulting subject matter experts. However, the growth of Participatory Modeling (PM) studies using techniques such as Fuzzy Cognitive Mapping or Causal Mapping [29, 36] has resulted in an abundance of systems maps across topic areas, which can thus be used by modelers to create a model operating within any subset of the map. To appreciate the coverage of systems maps, consider examples from health [23] such as the Foresight Obesity Map [19] (over 100 factors and 300 links) or our map on mental and physical well-being in relation to body weight [5] (269 links). Mid-sized maps are even more common, and may remain sufficient for many modeling applications, such as health technology adoption [24] (52 factors and 105 connections) or radiotherapy [27] (66 links). A plethora of maps has also been developed to study ecology and sustainability [20, 8, 15] or social challenges [18]. By clearly summarizing the salient constructs within the problem space, such maps present an opportunity to shift from an ad-hoc model building approach to a systematic and transparent one, thus decreasing structural uncertainty. Specifically, the use of a map and accompanying datasets allows to *systematically* answer three essential interrelated questions for model building:

¹ A structure of the problem space may also be called ‘domain model’ [33] or ‘conceptual model’. To avoid confusion on the multiple uses of ‘model’, we reserve this term for the *simulation* model, that is, the operational model that was obtained after making design decisions within the problem space.

- Which factors should we include?
- Where in the problem space is it necessary to make *assumptions*?
- What is the *resulting uncertainty* of the model?

In this paper, we propose a framework to express the task of creating a model as a matter of navigating a systems maps given a set of datasets. Using this framework, the strategies used by modelers to answer the three questions above are made both transparent and systematic using graph algorithms. Through this framework, we note that creating the perfect model is an NP-Hard problem, hence there is no perfect strategy to automatically generate a model: a heuristic needs to be chosen by modelers based on measures that they explicitly wish to favor. Our contributions are as follows:

- (1) We present the first mathematical framework to express model design choices based on data availability and their effects on uncertainty.
- (2) We explain why a perfect model cannot be automatically created (NP-Hard problem), hence emphasizing the need for heuristics.
- (3) We demonstrate how common heuristics used by modelers can be formulated in this framework, using a guiding example from a model of suicide.

The remainder of this paper is organized as follows. In section 3, we introduce our framework formally and exemplify its elements using a model for suicide prevention. In section 4, we explain why the model creation problem is NP-Hard and demonstrate the effects of common heuristics. Finally, we discuss the potential of shifting the process of model building from an ad-hoc unspecified approach to the use of transparent heuristics within a systems map.

2 Framework

Intuitively, modelers have access to a systems map describing the problem space as well as at least one dataset. For the model to be adequate, end users need to see the effects of simulated interventions, which requires each intervention to eventually impact at least one observable outcome through a chain of rules. The task of creating a model thus requires maintaining paths from interventions to observable outcomes through the problem space. Modelers use a *strategy* regarding which paths to take, in part based on data. Paths may travel through nodes for which modelers do not have data, thus they would need to make assumptions about the simulated entities' feature values for these nodes. Paths will also go through edges, which represent causal mechanisms that would be turned into simulation rules. Some of these edges may be intuitively understood (e.g., 'more trauma leads to more suicide ideation') but not sufficiently characterized to write model rules, thus leading to additional assumptions (e.g., every unit of trauma leads to an increase p in suicide ideation). A modeling strategy is thus an algorithm that identifies a subset of a systems map based on available data and makes assumptions to address unknown nodes and edges.

The main elements of the framework are listed in Table 1 and will now be explained along with their formal notation. The problem space is represented by

Table 1. Main elements of the framework

Notation	Meaning
$\mathbb{G} = (V, E)$	Problem space represented as a map \mathbb{G} , consisting of labeled nodes V and directed connections E
I	Interventions (or ‘inputs’) needed by model users
O	Observations (or ‘outputs’) that model users require to quantify the effects of their interventions
$w(e)$	Uncertainty value of a causal connection in the problem space
\mathbb{D}	Set of all data sources available to modelers
$Sources(v)$	Set of data sources that include the node v ; in other words, data sources that modelers can use to initialize v
S	A model design strategy. Given the map \mathbb{G} and data sources \mathbb{D} , it specifies which subset of the map to use in a model, which data source will be involved, and what level of uncertainty will arise. For a model to be adequate, a valid strategy must ensure that each intervention from I has an observable effect in O .

a directed, labelled, weighted causal graph $\mathbb{G} = (V, E)$. The nodes V represent factors in the problem space and potential candidates for the entities’ features, such as age, depression, or suicide ideation. To characterize the task of model building, we need to identify nodes that play particular roles. A subset of the nodes $I \subset V$ represents the intervention nodes (e.g., public health campaigns, economic interventions), also known as **inputs**. Similarly, the subset $O \subset V$ represents the **outcome** nodes (e.g., number of suicide attempts, prevalence of suicide ideation), also known as **outputs**. For a model to be viable, the interventions of interest need to eventually affect the outcomes; otherwise, the model does not offer support to examine the consequence of the intervention. For example, the evaluation of a suicide prevention package may measure the impact of economic interventions through a reduction in suicide attempts.

The edges E stand for causal connections and have an associated value (i.e., edge weight) denoted $w(e) \mapsto \mathbb{R}+, e \in E$. This value denotes the uncertainty associated with the edge. We encode the value as a positive real number rather than a boolean because systems maps *may* provide fine grained information on the *amount* of uncertainty, which can thus be expressed without needing to amend our framework. For example, Fuzzy Cognitive Maps can specify uncertainty [17] by measuring the extent to which there is disagreement about the causal strength of an edge (via entropy) among participants [12] or by checking in a corpus whether there is enough supporting evidence for each proposed connection [28, 22]. In the case of Fuzzy Grey Cognitive Maps, each edge has a Grey uncertainty [26]. In a coarse categorization, such as by asking subject-matter experts whether a factor impacts another [25], uncertainty would either be 1 (present; maximal) or 0 (absent; minimal).

Information on causal mechanisms is held at the level of the systems map via $w(e)$ because it represents a fact about the system, independently of any data source selected by modelers. This is reflective of the fact that a systems

map serves as a synthesis of the evidence base [19, 5, 24, 27]. For example, a map may stipulate that abusing or neglecting children has an impact on their risk for suicidal ideation, or that a suicide attempt can lead to death. Data sources may help to understand how these general mechanisms work in a specific population, but the mechanisms exist irrespective of a population. In contrast, information on the entities' features depend on the data sources used, thus allowing to capture how traits are expressed in specific populations. The collection of data sources available to modelers is denoted by $\mathbb{D} = \{D_1, \dots, D_n\}$. Each data source may hold information on some of the nodes. We denote the set of data sources for a node $v \in \mathbb{V}$ by $Sources(v)$. When $Sources(v) = \emptyset$, modelers have no data regarding this specific node hence its uncertainty is maximal and its inclusion in a model would come at the cost of making an assumption.

Modelers can employ one of several strategies \mathbb{S} to design a model. As shown in the next section, examples may include finding the simplest set of rules from each intervention to an outcome, or finding rules that avoid uncertainty whenever possible. A strategy $S \in \mathbb{S}$ is thus a function:

$$S : \underbrace{(\mathbb{G}, \mathbb{D})}_{\text{given map and datasets}} \mapsto \left(\underbrace{G = (V \subseteq \mathbb{V}, E \subseteq \mathbb{E})}_{\text{selects a map subset}}, \underbrace{D \subseteq \mathbb{D}}_{\text{dataset used}}, \underbrace{\mathbb{R}}_{\text{uncertainty cost}} \right)$$

In line with the earlier explanations in this section, a strategy is *valid* if and only if there is a path from every intervention node to at least one observable node. Formally, this condition is enforced by checking the following:

$$S \text{ is valid} \iff \forall i \in I, \exists o \in O \text{ s.t. } (i, v_1), \dots, (v_n, o) \in E$$

A core aim for modelers is to design a valid strategy while minimizing the uncertainty cost. The design of a model can thus be operationalized through this framework as a discrete optimization problem in a graph given a set of datasets.

3 The necessity and design of heuristics

3.1 The impossible quest for perfection in model design

Intuitively, an ideal model design is one that satisfies all needs of the end users while making the least number of assumptions. Formally, that would be an optimal strategy $S^* \in \mathbb{S}$ such that S^* is valid and its associated uncertainty cost is minimal among all valid strategies. However, finding the best strategy may not be feasible in practice as shown in the theorem below.

Theorem 1. *Identifying the best strategy S^* is an NP-Hard problem.*

Proof. To compute the optimal strategy S^* , we are given intervention and outcome nodes. The aim is to select nodes and edges in between (i.e., the 'inner part' of the network) such that the sum of selected nodes' and edges' weights is minimal, while maintaining connectivity. Minimizing this inner cost while providing connectivity is known as the *Minimum Spanning Tree with Inner nodes*

cost problem (MSTI). The MSTI problem is NP-Hard as a special case of the Connected Dominated Set problem [9, 21], hence our problem is also NP-Hard.

Alternatively, our problem can be related to the node weighted Steiner tree in which ‘terminals’ must be connected and the cost function is the sum of the nodes’ weights selected to provide this connectivity. Considering the special case in which there is a single outcome node and multiple intervention nodes, then they collectively form the ‘terminals’ used in a node weighted Steiner tree. Further consider that there is no uncertainty on any edge, hence leaving only the uncertainty on the nodes. Then, the special case is equivalent to minimizing a node weighted Steiner tree, which is an NP-Hard problem [4].

3.2 Imperfect yet practical: two common model heuristics

The previous section established that building model that is fit-for-purpose and minimizes uncertainty is an NP-Hard problem. However, modelers routinely build models, which implies that they use heuristics. In this section, we show how two such heuristics can be expressed in our framework, thus making the model building process more transparent and systematic. As example for both heuristics is provided in Figure 2.

One strategy employed by modelers is to ‘Keep It Simple Stupid’ (KISS) in which “one *only* tries a more complex model if simpler ones turn out to be inadequate” [6]. In other words, the model design is justified by ‘the sake of simplicity’. Since each intervention must be measured via an observation, a translation of KISS into a process would be to find the *simplest* way of connecting each intervention to one observation (possibly the same). When that strategy is expressed algorithmically, it is equivalent to using the *shortest path* from each intervention to one observation. Algorithm 1 formalizes this strategy by using a Breadth-First Search to generate each shortest path. We make two observations:

- (1) Focusing on the shortest number of rules will not take the locations of unknowns into consideration. A *slightly* more complex model in number of rules may have gone through a path better supported by data, thus producing a model that needs fewer assumptions.
- (2) *Independently* finding a path for each intervention can produce a model that is simple for each intervention, but overall much more complex than strictly necessary. For instance, taking a short detour for one intervention may have allowed it to share the rest of the journey toward an observation using the path of another intervention. Sharing paths or finding ‘synergies’ may thus help to keep the whole model simpler and potentially lower its uncertainty.

A second approach captures the preference of modelers who actively avoid creating poorly understood rules. Their preference is not to keep the model ‘simple’: rather, they seek a more robust model in which rules can be backed by evidence as much as possible. From an algorithmic standpoint, this consists of generating paths from every intervention node to the observation nodes and selecting the one with the least overall uncertainty. This selected path will thus

use a sequence of factors with well-understood relationships to lead from an intervention node to an observation node.

Algorithm 1: Generate all shortest paths using Breadth-First Search

Input: List of edges in causal map, List of intervention nodes, List of observation nodes

```

1 paths = {∅}, finalPaths = ∅ // create empty list of lists and a
  map
2 foreach intervention do
3   | paths ← paths ∪ {{intervention}}
4 while ∃intervention ∉ finalPaths // continue until we have a path
  for each intervention
5 do
6   | newPaths = {∅} // for each current path, look for next
  steps
7   foreach path  $i_1, \dots, i_n \in$  paths do
8     | targets ← { $t \mid (i_n, t) \in$  edges}
9     | foreach target ∈ targets // for each possible next step,
  save the potential path
10    | do
11    | | newPaths ← path ∪ target
12    | foreach path  $i_1, \dots, i_n \in$  newPaths do
13    | | if  $i_1 \notin$  finalPaths // ignore path if we already have a
  finalPath for its root
14    | | then
15    | | | if  $i_n \in$  observations // if the path is complete, save
  to finalPaths
16    | | | then
17    | | | | finalPaths( $i_1$ ) ← path
18    | | | else
19    | | | | paths ← paths ∪ path // otherwise keep path for
  next iteration

```

This strategy may lead to long and awkward tangents. For example, instead of creating few rules with limited empirical or theoretical backing (e.g., after-school programs reduce school problems thus reducing suicide ideation), this strategy may result in creating *many* rules: after-school programs reduce school problems thus reducing parental frustration, which means parents may be less likely to cope with frustration through substance use, hence children are less exposed to unhealthy coping strategies, thus they can deal better with their own issues and overall are less likely to engage in suicidal thoughts. Each of these rules may be backed by stronger evidence, but the collective chain of rules produces a meandering model that can appear less plausible overall. As modelers may avoid *both* arbitrarily long chains of rules *and* models whose rules lack evidence, a strategy would need to discourage both uncertainty and long paths. Algorithm 2 implements this approach by assigning a unit cost of 1 to all known edges and a penalty value to unknown edges.

Algorithm 2: Generate all minimal paths using Dijkstra's algorithm

Input: List of edges in causal map, List of intervention nodes, List of observation nodes, Penalty

```

1 finalPaths  $\leftarrow \emptyset$ 
2 foreach intervention do
3     visitedNodes  $\leftarrow \emptyset$ 
4     paths  $\leftarrow \{\{intervention, 0\}\}$  // each node in a path is a pair
        (name, cost)
5     bestTarget  $\leftarrow \{0, 0\}$ 
6     while  $bestTarget_1 \notin observations$  do
7         bestTarget  $\leftarrow \{\emptyset, \infty\}$  // track current most-optimal edge
            to add
8         bestPath  $\leftarrow \{\emptyset\}$ 
9         foreach  $path \in paths$  do
10            traversal  $\leftarrow \emptyset$  // track current traversal through path
11            foreach  $node \in path$  do
12                traversal  $\leftarrow traversal \cup node$ 
13                targets  $\leftarrow \{t | (node_1, t) \in edges\}$ 
14                foreach  $target \in targets$  do
15                    if  $target \notin visitedNodes$  // ignore
                        already-visited nodes
16                    then
17                        if  $(node_1, target)_{weight} \in ]-1, 1[$  then
18                            cost  $\leftarrow node_2 + 1.0$  // if edge weight is
                                known, marginal cost = 1
19                        else
20                            cost  $\leftarrow node_2 + penalty$  // otherwise,
                                marginal cost = penalty
21                        if  $cost < bestTarget_2$  then
22                            bestTarget  $\leftarrow target$ 
23                            bestPath  $\leftarrow traversal$ 
24            newPath  $\leftarrow bestPath \cup bestTarget$ 
25            if  $bestTarget \in observations$  // if target  $\in$  observations,
                this root is finished
26            then
27                finalPaths(intervention)  $\leftarrow newPath$ 
28            else
29                paths  $\leftarrow paths \cup newPath$  // otherwise, save to
                    current paths and continue
30                visitedNodes  $\leftarrow visitedNodes \cup bestTarget$ 
    
```

Both algorithms are exemplified on the next page. They produce the same paths for interventions $I1$ and $I2$ but differ on $I3$. In this example, Algorithm 2 is better: it uses 10 evidence-based rules and makes five assumptions, whereas Algorithm 1 needs more rules (11) and also makes more assumptions (six).

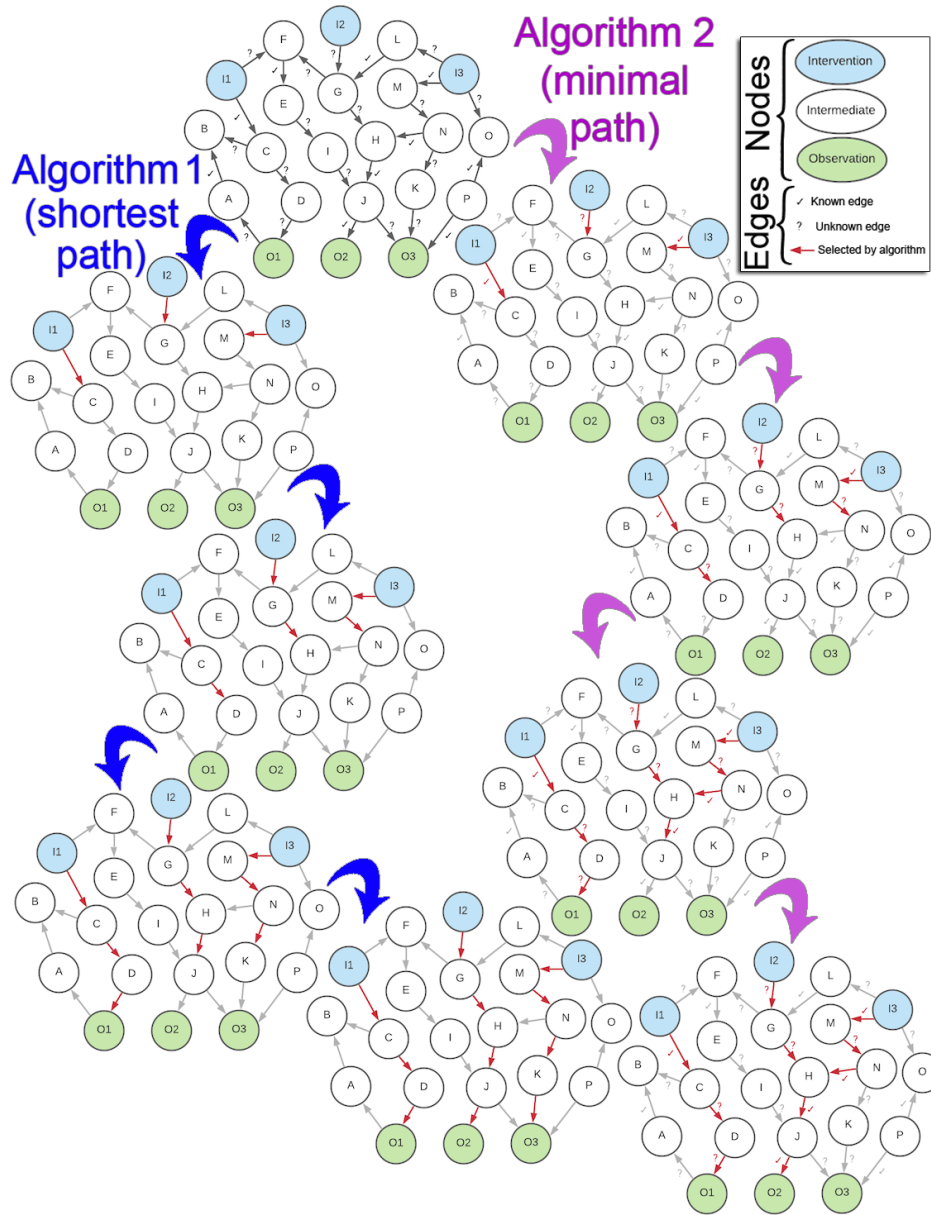


Fig. 2. The same map (top) being processed by both algorithms: shortest paths (left) and minimal paths (right). *This high resolution figure can be zoomed in using a digital copy of this article.*

3.3 Differences in practice: a sample case study

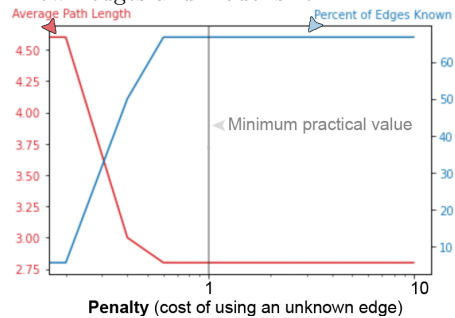
Algorithms 1 and 2 are *conceptually* different as they correspond to different preferences from modelers: the simplest set of rules to connect each intervention (Alg. 1) or an emphasis on evidence-based rules (Alg. 2). As exemplified by Figure 2 on the previous page, these algorithms can produce different models in terms of number of rules or assumptions. Two items remain in order to assess the impact of these different model design strategies. First, we need to examine how the number of rules or assumptions differ when these strategies are used *in practice* rather than in an idealized environment. Second, we need to evaluate the possible *choices for datasets* regarding the concept nodes, which also benefits from a case study rather than an artificially constructed situation.

Our brief case study is a model of suicide. Data is openly accessible at <https://osf.io/7npx4/>, which include (i) the systems map together with the uncertainty level for each edge (3-‘big map with weights’) and (ii) which one(s) of four datasets can be used for each concept node (5-‘Nodes data availability’). Our algorithms are implemented in a Jupyter Notebook for Python 3, which also contains the results (6-‘Common modeling strategies’). The systems map is composed of 361 concept nodes and 946 causal edges. As detailed in the Notebook, there are 5 intervention nodes and 3 outcome nodes.

We performed a parameter sweep to assess how the penalty in Algorithm 2 impacts its two target outcomes: the percentage of unknown edges and the average size of the model. Results in Figure 3 show that the only difference is encountered if the cost of unknowns is less than 1, that is, less than using a *known* edge. Specifically, the only values of penalty that had an effect on results were in the range $0.2 \leq \text{penalty} \leq 0.6$. However, creating evidence-based rules cannot be more expensive than making assumptions for unknown edges. The penalty of unknown would thus be always greater than 1. Consequently, when the penalty is within a *valid range*, then the results are always the same.

Using Algorithm 1, the model requires 3 assumptions on edges and 4 assumptions on nodes. The five paths created for each intervention had no overlap. Two datasets can be interchangeably used as they each support the same number of nodes. With Algorithm 2, *results are identical*. Despite valuing different aspects of a model and having a large problem space, the two design strategies result in the same model. Either strategy has room for improvement as the resulting model requires making many assumptions (7) and is composed of parts that may be locally optimal but miss an opportunity for global savings (disjoint paths across all interventions).

Fig. 3. Evaluation of the penalty parameter in Algorithm 2 with respect to unknown edges and model size.



4 Discussion and conclusion

Several frameworks such as the Characterisation and Parameterisation (CAP) framework [31] and its revised version [30] have been proposed to track how computational individual-level models are designed. Such frameworks provide valuable guidance to help standardize the description of modeling studies (c.f. the application of the CAP framework to 11 studies in [32]) by listing the broad types of data or methods involved at each stage. However, there is currently no framework to explain how modelers choose what factors to keep and which rules to make, based on the goal of the model and available data. These choices are important for the transparency and replicability of modeling studies, which often face a high level of structural uncertainty. In addition, these choices have consequences on the robustness and computational costs of models: inefficient model building strategies that result in making a very large number of assumptions may require extensive sensitivity analyses, which come at a high cost. In this paper, we propose the first formal framework to express model-building strategies. We demonstrated that the perfect strategy is NP-Hard, hence modelers will employ heuristics that emphasize specific aspects. We stress this point, as it means that a perfect model cannot be built automatically and instead modelers need to clearly state which measures they value, then map these preferences onto an algorithm. We showed how two common strategies can be systematically expressed in an algorithmic manner and demonstrated that they *can* result in different models, although differences may not be manifest *in practice*.

We noted that both of these common model-building strategies miss several opportunities to decrease the uncertainty of the resulting model. Creating the model with the least number of rules to keep it ‘simple’ may result in high uncertainty compared to having a tolerance for slightly more rules as long as they are evidence-based. Most importantly, when models are designed to support multiple interventions, it would be beneficial to aim for synergies by identifying common mechanisms across interventions rather than operationalizing each one independently. The design and evaluation of algorithms leveraging these opportunities would be a worthwhile investment for future studies. Such tools can help modelers in shifting from the currently time consuming and ad-hoc practice of model design into a more efficient and systematic approach.

As our framework is the first to tackle complex practices, it comes with simplifications. When simplifications prevent a direct use of the framework by a modeling team, they become limitations and changes are necessary. Although we separately report uncertainty on causal mechanisms (i.e. on edges in the problem space) and concept nodes (i.e. insufficient data), our framework is limiting in capturing the cost of uncertainty on nodes. We focused on using *one* dataset to maximize node coverage, but modelers can use *multiple* datasets as long as the individuals captured in these datasets can be accurately linked. For example, consider that two datasets have nothing in common: one cover depression and stress, while the other contains information on poverty and bullying. If a simulated entity was independently given baseline values on depression and poverty, then it would ignore the correlation between these features and model uncer-

tainty grows. Conversely, if the two datasets have shared features that strongly characterize individuals (e.g., age, gender, income, ethnicity) then we may preserve more (but not all) of the real-world dependencies, hence lowering model uncertainty. A possible extension of our framework would thus address how the cost of uncertainty is impacted by the ability to link datasets.

References

1. Abdou, M., Hamill, L., Gilbert, N.: Designing and Building an Agent-Based Model, pp. 141–165. Springer Netherlands, Dordrecht (2012)
2. Bojke, L., et al.: Characterizing structural uncertainty in decision analytic models: a review and application of methods. *Value in Health* **12**(5), 739–749 (2009)
3. Carley, K.M.: Social-behavioral simulation: Key challenges. *Social-Behavioral Modeling for Complex Systems* (2019)
4. Demaine, E.D., Hajiaghayi, M., Klein, P.N.: Node-weighted steiner tree and group steiner tree in planar graphs. *ACM Transactions on Algorithms (TALG)* **10**(3), 1–20 (2014)
5. Drasic, L., Giabbanelli, P.J.: Exploring the interactions between physical well-being, and obesity. *Canadian Journal of Diabetes* **39**, S12–S13 (2015)
6. Edmonds, B., Moss, S.: From kiss to kids—an ‘anti-simplistic’ modelling approach. In: *International workshop on multi-agent systems and agent-based simulation*. pp. 130–144. Springer (2004)
7. Epstein, J.M.: Why model? *Journal of Artificial Societies and Social Simulation* **11**(4), 12 (2008)
8. Firmansyah, H.S., et al.: Identifying the components and interrelationships of smart cities in indonesia: Supporting policymaking via fuzzy cognitive systems. *IEEE Access* **7**, 46136–46151 (2019)
9. Fleischer, R., et al.: Approximating spanning trees with inner nodes cost. In: *Sixth International Conference on Parallel and Distributed Computing Applications and Technologies (PDCAT’05)*. pp. 660–664. IEEE (2005)
10. Giabbanelli, P.J., Jackson, P.J.: Using visual analytics to support the integration of expert knowledge in the design of medical models and simulations. *Procedia Computer Science* **51**, 755–764 (2015)
11. Giabbanelli, P.J., et al.: Ideal, best, and emerging practices in creating artificial societies. In: *2019 Spring Simulation Conf. (SpringSim)*. pp. 1–12. IEEE (2019)
12. Giabbanelli, P., Torsney-Weir, T., Mago, V.: A fuzzy cognitive map of the psychosocial determinants of obesity. *Applied Soft Computing* **12**(12), 3711–3724 (2012)
13. Gilbert, N., Doran, J.: *Simulating societies: the computer simulation of social phenomena*. Routledge (2018)
14. Gilbert, N., Stoneman, P.: *Researching social life*. Sage (2015)
15. Gray, S., et al.: Assessing (social-ecological) systems thinking by evaluating cognitive maps. *Sustainability* **11**(20), 5753 (2019)
16. Grimm, V., et al.: The odd protocol for describing agent-based and other simulation models: A second update to improve clarity, replication, and structural realism. *Journal of Artificial Societies and Social Simulation* **23**(2) (2020)
17. Lavin, E.A., Giabbanelli, P.J.: Analyzing and simplifying model uncertainty in fuzzy cognitive maps. In: *2017 Winter Simulation Conference (WSC)*. pp. 1868–1879. IEEE (2017)

18. Mago, V.K., et al.: Analyzing the impact of social factors on homelessness: a fuzzy cognitive map approach. *BMC medical informatics and decision making* **13**(1), 94 (2013)
19. McPherson, K., Marsh, T., Brown, M.: Foresight report on obesity. *The Lancet* **370**(9601), 1755 (2007)
20. Mourhir, A.: Scoping review of the potentials of fuzzy cognitive maps as a modeling approach for integrated environmental assessment and management. *Environmental Modelling & Software* p. 104891 (2020)
21. Peng, C., Tan, Y., Zhu, H.: On computing the backbone tree in large networks. In: 2008 IEEE Systems and Information Engineering Design Symposium. pp. 118–122. IEEE (2008)
22. Pillutla, V.S., Giabbanelli, P.J.: Iterative generation of insight from text collections through mutually reinforcing visualizations and fuzzy cognitive maps. *Applied Soft Computing* **76**, 459–472 (2019)
23. de Pinho, H.: Generation of systems maps: Mapping complex systems of population health. *Systems Science and Population Health* pp. 61–76 (2017)
24. Rahimi, N., et al.: Soft data analytics with fuzzy cognitive maps: modeling health technology adoption by elderly women. In: *Advanced Data Analytics in Health*, pp. 59–74. Springer (2018)
25. Reddy, T., Giabbanelli, P.J., Mago, V.K.: The artificial facilitator: guiding participants in developing causal maps using voice-activated technologies. In: *International Conference on Human-Computer Interaction*. pp. 111–129. Springer (2019)
26. Salmeron, J.L.: Modelling grey uncertainty with fuzzy grey cognitive maps. *Expert Systems with Applications* **37**(12), 7581–7588 (2010)
27. Salmeron, J.L., Papageorgiou, E.I.: A fuzzy grey cognitive maps-based decision support system for radiotherapy treatment planning. *Knowledge-Based Systems* **30**, 151–160 (2012)
28. Sandhu, M., Giabbanelli, P.J., Mago, V.K.: From social media to expert reports: The impact of source selection on automatically validating complex conceptual models of obesity. In: *International Conference on Human-Computer Interaction*. pp. 434–452. Springer (2019)
29. Schmitt-Olabisi, L., McNall, M., Porter, W., Zhao, J.: *Innovations in Collaborative Modeling*. MSU Press (2020)
30. Smajgl, A., Barreteau, O.: Framing options for characterising and parameterising human agents in empirical abm. *Environmental Modelling & Software* **93**, 29–41 (2017)
31. Smajgl, A., et al.: Empirical characterisation of agent behaviours in socio-ecological systems. *Environmental Modelling & Software* **26**(7), 837–844 (2011)
32. Smajgl, A., Barreteau, O.: *Empirical agent-based modelling-challenges and solutions*, vol. 1. Springer (2014)
33. Stepney, S., Polack, F.A.: Discovery phase patterns: building the domain model. In: *Engineering Simulations as Scientific Instruments: A Pattern Language*, pp. 117–145. Springer (2018)
34. Swarup, S.: Adequacy: what makes a simulation good enough? In: 2019 Spring Simulation Conference (SpringSim). pp. 1–12. IEEE (2019)
35. Treasury, H.: *The aqua book: guidance on producing quality analysis for government* (2015)
36. Voinov, A., et al.: Tools and methods in participatory modeling: Selecting the right tool for the job. *Environ Modell Softw* **109**, 232–255 (2018)