# GRAPH-THEORETIC AUTOFILL

#### MICHAEL MAYER AND DOMINIC VAN DER ZYPEN

ABSTRACT. Imagine a website that asks the user to fill in a web form and – based on the input values – derives a relevant figure, for instance an expected salary, a medical diagnosis, or the market value of a house. How to deal with missing input values at run-time? Besides using fixed defaults, a more sophisticated approach is to use predefined dependencies (logical or correlational) between different fields to autofill missing values in an iterative way. Directed loopless graphs (in which cycles are allowed) are the ideal mathematical model to formalize these dependencies. We present two new graph-theoretic approaches to filling missing values at run-time.

#### 1. Introduction

The Internet offers many online calculators that provide relevant figures based on the values entered in a web form. Examples are salary calculators, web-based medical advice, and tools to compute the typical rent of an appartment, just to name a few. Usually, all input fields are mandatory, i. e. cannot be left blank by the user. While this minimizes programming effort by the website provider, it might force the user to make up or guess unknown information (like the living area of a 4-room appartment whose monthly rent is to be estimated by the calculator) just for the sake of completeness, or the user will even search the web for a simpler calculator without ever returning.

More user-friendly web forms offer fixed default values at least for some of the fields. On one hand, this approach is attractive thanks to its low programming effort. On the other hand, the stronger the input fields are interrelated, the less appropriate a fixed default might be in certain cases. While a default living area of  $80\,m^2$  might serve as a good default for a typical appartment, it is certainly unrealistic for a studio or a 10-room penthouse with indoor swimming-pool.

What are alternatives to handle missing input values on web forms? One option is the reduced feature approach from statistical predictive modelling, see e. g. [5] or [6]. There, for each combination of available fields, an individual calculator is applied. The drawback is obvious. Even for a low number p of optional input fields, the programming effort exploses, as  $2^p$  calculators have to be derived, implemented and supported.

An elegant compromise between offering over-simplistic fixed defaults and the unfeasible reduced feature approach is to autofill missing values by prespecified functions of other input values, i. e. using dynamic defaults for certain fields. This can either be made visible to the user or invisibly applied before calling the underlying formula of the calculator.

Consider as example a calculator for the ideal weight based on body height, sex and age. There, the autofill strategy could be as follows.

- Sex s: Mandatory (1: male, 0: female)
- Age x: Autofilled by height z in cm using the formula

$$x = f(z) := \begin{cases} \lfloor (z - 30)/130 \cdot 16 + 1 \rfloor & \text{if } 30 \le z \le 160, \\ 40 & \text{if } z > 160, \\ 1 & \text{if } z < 30. \end{cases}$$

• Body height z: Autofilled by height and sex by

$$z = g(x,s) := \begin{cases} 162 + 16s & \text{if } x > 16, \\ \lfloor (x-1)/16 \cdot 130 + 30.5 \rfloor & \text{if } x \le 16. \end{cases}$$

Thus, the web form with the fields "age", "sex" and "height" can be considered to be filled (in the sense as the underlying weight formula can be applied) as soon as sex and either of height and age is provided. Or to express the same thought in its negative sense: Even if we have specified replacement functions for age and height, the web form cannot be autofilled if age as well as height are left blank. In more complex situations, also the order in which the replacement functions are applied could matter, e.g. if sex would be autofilled by a function of the height.

The purpose of this article is to use the notion of directed graphs (in which cycles are allowed) to provide a theoretical framework to determine if an arbitrarily complex web form (or any other multivariate input, e.g. the argument list of a function written in a programming language like C) can be filled by a fixed set of replacement functions and the partial input provided so far by the user. To do so, each input field is associated with a vertex in a graph and each replacement function with at least one directed edge between these vertices. Note that our considerations do not depend on how the replacement functions are defined or how accurate they are. In practice they will be chosen based on expert knowledge, literature, logical rules or by exploring statistical relationships.

Suppose we are given an observation in which knowing the value of field A would enable you to guess the value of field B. A natural way to represent this is to draw an **arrow** from A to B:

$$A \to B$$

Notice that the relation "B can be guessed if we know A" is sometimes asymmetric, that is, only works in one direction, as the following example illustrates:

Let A stand for "gender" and B for "number of pregnancies had so far". Then if you know that some observation represents a male then you can infer that the value of B must be 0. On the other hand if we know that the number of pregnancies is 0, the individual at hand can be either male or female, depending on the dataset at hand even with roughly equal probability.

So this shows that it is natural to choose *directed graphs*, essentially points connected by arrows, as a model for representing the conclusions that can be made between some fields based on replacement functions. The terms of a directed graph, and other terms, will be defined rigorously in the next section, and for a good introduction into graph theory, we point the reader to [1].

The left side of Figure 1 shows the graph corresponding to the example of the weight-calculator.

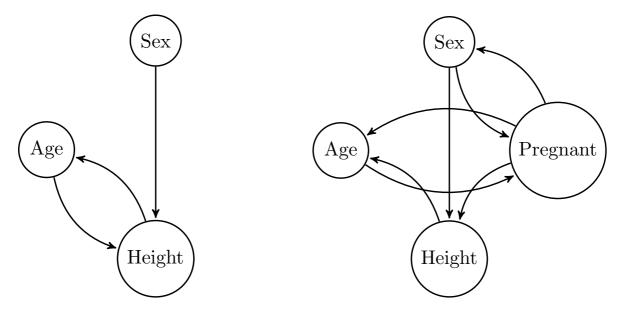


FIGURE 1. On the left: Graph for weight-calculator. On the right: A more complex situation with additional field.

We consider two kinds of replacement functions. In Section 2 we look at the case of "complete determination", that is we are only allowed to calculate the value of a certain field Y if all the values of the fields from which an arrow leads into Y are known or already derived from other fields. Thus, a replacement function can only be evaluated if all arguments are made available in some way as in the case of the weight-calculator. Next, Section 3 applies to the situation of "partial determination", where we are allowed to calculate the value of a certain field Y if only some of the values of the fields from which an arrow leads into Y are known or derived from other fields, i. e. some missing or undetermined arguments in the replacement functions are allowed.

# 2. FILLING WITH COMPLETE DETERMINATION

A directed self-loopless graph (which we will refer to as directed graph for simplicity) is a tuple G = (V, E) where V is a finite non-empty set and

$$E \subseteq V \times V \setminus \{(v, v) : v \in V\}.$$

We call the set V the set of *vertices*, and they represent the data fields or variables. An *edge* represents a determination arrow: If  $(v, w) \in V$  we can make some conclusion from the variable v to variable w.

**Definition 2.1.** Let  $x, y \in V$ . Then there is a directed path from x to y if there is  $n \in \mathbb{N}$  and a map  $p : \{0, \ldots, n\} \to V$  such that

- (1) p(0) = x, p(n) = y;
- (2) for  $k \in \{0, ..., n-1\}$  we have  $(p(k), p(k+1)) \in E$ .

If  $n \ge 1$  and  $p: \{0, \ldots, n\} \to V$  is a directed path from x to x we call  $p(\{0, \ldots, n\}) \subseteq V$ , that is the image of p, a cycle.

Note that the above definition implies that a cycle has at least 2 elements.

**Definition 2.2.** For  $v \in V$  we denote the set of incoming vertices by  $In(v) = \{w \in V : (w,v) \in E\}$  and we let  $A(G) = \{v \in V : In(v) = \emptyset\}$ .

A(G) represents the fields for which no replacement function are defined. Thus, they are always mandatory (or, without affecting our theory, need a constant default).

**Definition 2.3.** Let  $I \subseteq V$  and let  $v \in V$ . Then we say that I determines the vertex v if  $In(v) \neq \emptyset$  and  $In(v) \subseteq I$ .

In other words, I stands for the arguments of the replacement function for v which is represented in the graph by the determination arrows that come in from In(v).

**Definition 2.4.** We denote the set of vertices determined by I by Dtm(I). Inductively we define the "determination closure" of  $I \subseteq V$  in the following way:

- (1)  $I_0 := I$ ;
- (2) for  $n \geq 0$  set  $I_{n+1} := I_n \cup \operatorname{Dtm}(I_n)$ .

We say that  $I \subseteq V$  fills V if there is  $n \in \mathbb{N}$  such that  $I_n = V$ . Note that trivially, V itself fills V.

Consequently, if all fields associated with I are entered, then the remaining fields can be autofilled iteratively by the prespecified replacement functions.

The following lemmata are mainly used to derive a characterization theorem for filling subsets.

**Lemma 2.5.** If  $I \subseteq V$  fills V then  $A(G) \subseteq I$ .

*Proof.* Trivial, follows from definition of determination.

**Lemma 2.6.** Suppose  $I \subseteq V$  and suppose that C is a cycle with  $I \cap C = \emptyset$ . Then

$$I_1 \cap C = (I \cup Dtm(I)) \cap C = \emptyset.$$

*Proof.* We prove the contrapositive. Suppose there is  $c^* \in C$  such that  $c^* \in I_1$ . Since C is a cycle, we have

$$C \cap \operatorname{In}(c^*) \neq \emptyset$$
,

say  $d \in C \cap \operatorname{In}(c^*)$ . Now  $c^* \in I_1$  implies  $\operatorname{In}(c^*) \subseteq I$  by definition, therefore

$$d \in \operatorname{In}(c^*) \cap C \subseteq I \cap C$$
,

so 
$$I \cap C \neq \emptyset$$
.

An inductive application of Lemma 2.6 shows that any subset that fills V intersects every cycle in G = (V, E):

**Lemma 2.7.** If  $I \subseteq V$  fills V and  $C \subseteq V$  is a cycle, then  $I \cap C \neq \emptyset$ .

*Proof.* Let  $I\subseteq V$  be any subset of V and let  $C\subseteq V$  be a cycle. Applying Lemma 2.6 inductively implies that

 $(\star)$  if  $I \cap C = \emptyset$  then  $I_n \cap C = \emptyset$  for all  $n \in \mathbb{N}$ .

So if I fills V then  $I_n = V$  for some n, and in particular  $I_n \cap C \neq \emptyset$  for that n. The contrapositive of  $(\star)$  directly implies  $I \cap C \neq \emptyset$ .

This helps us to prove a characterization theorem for filling subsets  $I \subseteq V$ .

**Theorem 2.8.** Let G = (V, E) be a self-loopless directed graph, and let  $I \subseteq V$ . Then the following statements are equivalent:

- (1) I is filling;
- (2)  $A(G) \subseteq I$ , and for every cycle C in (V, E) we have  $C \cap I \neq \emptyset$ .

A web form can thus be autofilled as soon as values are provided

- (1) for all fields without replacement function, and
- (2) for one field per cycle.

Proof of Theorem 2.8. (1)  $\implies$  (2). Taken care of by Lemmas 2.5 and 2.7.

(2)  $\Longrightarrow$  (1). We assume that  $I \subseteq V$  is not filling and  $v \in I$  for every  $v \in V$  with  $\text{In}(v) = \emptyset$  and construct a cycle  $C^*$  that does not intersect I (i.e.  $C^* \cap I = \emptyset$ ).

Since I is not filling, we have  $I_n \neq V$  for all  $n \in \mathbb{N}$ . Since V is finite, the increasing sequence

$$I = I_0 \subseteq I_1 \subseteq I_2 \subseteq \dots$$

stabilizes at some  $N \in \mathbb{N}$ , that is there is  $N \in \mathbb{N}$  such that  $I_k = I_N$  for all  $k \geq N$ , and of course  $I_N \neq V$ .

So we pick  $z_0 \in V \setminus I_N$ .

Note that  $\operatorname{In}(z_0) \neq \emptyset$  because all vertices v with  $\operatorname{In}(v) = \emptyset$  are included in I by assumption, and  $z_0 \notin I$ . Since  $I_{N+1} = I_N$  we have that  $z_0$  is not determined by  $I_N$  which means  $\operatorname{In}(z_0) \not\subseteq I_N$ . So we pick  $z_1 \in \operatorname{In}(z_0) \setminus I_N$ .

Inductively, and similarly to above, we pick  $z_{k+1} \in \text{In}(z_k) \setminus I_N$  for all  $k \in \mathbb{N}$ .

Then we consider the set  $Z = \{z_k : k \in \mathbb{N}\}$ . By construction Z has empty intersection with  $I_N$  and therefore also with  $I \subseteq I_N$ . Moreover we have  $(z_{k+1}, z_k) \in E$  for all k. Because V is finite, the set Z must contain a cycle  $C^*$  and we have  $C^* \cap I = \emptyset$ .

A cycle C is said to be *minimal* if it is minimal amongst all cycles with respect to  $\subseteq$  (that is, whenever  $C' \subseteq C$  and C' is a cycle, then C' = C). Sometimes, minimal cycles are called *elementary*, for instance in [2].

**Remark 2.9.** It is easy to see that Theorem 2.8 can be made a bit simpler: in order to verify that a certain subset  $I \subseteq V$  is filling, it suffices to check that  $A(G) \subseteq I$  and that I intersects every minimal cycle.

**Example 2.10.** In the introductory example of a weight-calculator, we have considered the three input fields "Sex", "Age" and "Height" along with two replacement functions. As shown on the left hand side of Figure 1, the fields are represented by the three vertices in the graph G. The directed edge from "Height" to "Age" stands for the replacement function f for age based on height. Finally, the replacement function g for height depending on age and sex

is represented by the two other edges pointing to "Height". By our definitions, we can for instance make the following statements:

- (1) The vertex set {Sex, Age} determines {Height}.
- (2) The vertex {Height} determines {Age}.
- (3) The set A(G) equals  $\{Sex\}$  (no edge points to it).
- (4) By Theorem 2.8, the vertex set {Sex} is not filling, i. e. by entering only sex, the other fields cannot be autofilled. The reason is that the set {Sex} does not intersect the cycle formed by the two vertices {Age, Height}.
- (5) There are exactly three filling subsets for G: {Sex, Age}, {Sex, Height} and (trivially) also {Sex, Age, Height} as they all contain A(G) and at least one element of the only cycle.

For complex graphs, the application of Theorem 2.8 might need algorithmic support to identify A(G) and particularly the cycles to check during runtime if partial input already fills the remaining fields. To do so, the following algorithms can be applied.

**Algorithm 2.11** (Identifying A(G)). A convenient representation for a directed graph is the adjacency matrix: The vertex are numbered  $1, \ldots, n$  and we assign a binary  $n \times n$ -matrix

$$M_G \in \mathbb{Z}_2^{n \times n}$$

to G by setting  $M_G[i,j]=1$  if and only if  $(i,j)\in E$ , and  $M_G[i,j]=0$  otherwise.

Now A(G) is easily identified:  $i \in A(G)$  if and only if  $M_G[\cdot, i]$  (that is the *i*th column vector) is the constant 0 vector.

**Algorithm 2.12** (Identifying cycles). Different algorithms exist for finding cycles in a graph, see e. g. [2] for a solution and further references.

The next result link graph filling to the concept of directed acyclic graph (DAG).

**Theorem 2.13** (Connection to DAGs). Let G = (V, E) be a self-loopless directed graph, and let  $I \subseteq V$ . Furthermore denote by G' = (V, E') with  $E' = E \setminus \{(x, i) : x \in V \text{ and } i \in I\}$  the subgraph without edges pointing to any  $v \in I$ . Then the following are equivalent:

- (1) I is filling in G = (V, E);
- (2) I is filling in G' = (V, E') and G' is a DAG.

Proof. (1)  $\Longrightarrow$  (2). First, we prove that I is filling in G = (V, E'). Let  $I'_n$  be the determination closures of I in the graph G'. With  $I_n$  we denote the determination closures of I in G. For  $v \in V$  we denote by  $\operatorname{In}'(v)$  the set of incoming vertices in the graph G', and by  $\operatorname{In}(v)$  the set of incoming vertices in G. As I is filling in G by assumption, we have  $I_N = V$  for some  $N \in \mathbb{N}$ , so it suffices to show that

$$I'_n \supseteq I_n$$
 for all  $n \in \mathbb{N}$ .

We proceed by induction. Clearly  $I'_0 = I_0 = I$ . Suppose we have  $I'_k \supseteq I_k$  for some k and show that  $I'_{k+1} \supseteq I_{k+1}$ . Let  $v \in I_{k+1}$ . If  $v \in I_k$  we get  $v \in I'_{k+1}$  automatically since  $v \in I_k \subseteq I'_k \subseteq I'_{k+1}$ . If  $v \notin I_k$  then we have trivially  $v \notin I$ . Moreover, the fact that  $I_k$  determines v in G implies in particular that  $\operatorname{In}(v) \neq \emptyset$ , so we pick some incoming vertex  $j \in \operatorname{In}(v)$ . Since  $v \notin I$  we have  $(j,v) \in E'$  by the definition of E', so

(A) 
$$\operatorname{In}'(v) \neq \emptyset \text{ in } G'.$$

Since  $I_k$  determines v in the graph G, we get  $\operatorname{In}(v) \subseteq I_k$ , and by definition of G' we get  $\operatorname{In}'(v) \subseteq \operatorname{In}(v)$ . Combining this gives  $\operatorname{In}'(v) \subseteq I_k \subseteq I'_k$  by induction assumption, and together with (A) this implies that  $I'_k$  determines v in the graph G', so  $v \in I'_{k+1}$ , which finishes the proof of (1).

Next, we show that G' is a DAG. Assume that G' contains a cycle C. By Theorem 2.8 we know that  $I \cap C \neq \emptyset$ , so pick  $v^* \in I \cap C$ . Because C is a cycle there is a bijection  $p : \{0, \ldots, n\} \to C$  such that  $(p(k), p(k+1)) \in E'$  for  $k \in \{0, \ldots, n-1\}$  and  $(p(n), p(0)) \in E'$ . We can pick p such that  $p(0) = v^*$ . But by definition of E' we have  $(p(n), v^*) \in E \setminus E'$ , contradiction.

(2)  $\Longrightarrow$  (1). This is easily verified by noticing that Dtm(I) in the graph G' = (V, E') equals Dtm(I) in G = (V, E) (we don't even need the assumption that G' = (V, E) is a DAG).  $\square$ 

Since a DAG G' = (V, E') does not contain cycles, Theorem 2.8 ensures that  $I \subseteq V$  is filling if and only if  $A(G') \subseteq I$ .

This provides a second possibility to verify if a subset  $I \subseteq V$  of vertices of a directed graph G = (V, E) is filling or not: Delete all edges pointing to vertices in I and check if the resulting subgraph G' is a DAG with  $A(G') \subseteq I$ . If yes, I is filling.

**Example 2.14.** In Example 2.10 we have utilized Theorem 2.8 to show that, amongst others, the set  $I := \{\text{Sex}, \text{Age}\}$  is filling. Alternatively, we consider the subgraph G' without edges pointing to I (see left side of Figure 2.14). Obviously, G' is a DAG with I = A(G') and thus, I is filling in the original graph G. The right side of Figure 2.14 shows that  $I' = \{\text{height}\}$  alone is not filling: Although the subgraph G'' is a DAG, the subset  $A(G'') = \{\text{Sex}, \text{Height}\} \not\subseteq I'$ .

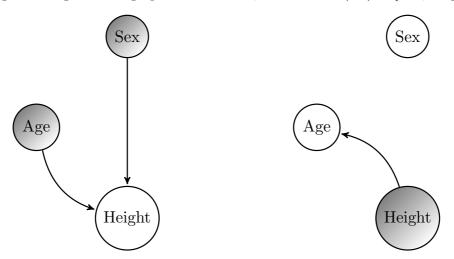


FIGURE 2. On the left: Subgraph G' without edges pointing to candidate set  $I := \{\text{Sex}, \text{Age}\}$  (highlighted). Right: Subgraph G'' without edges pointing to  $I' = \{\text{Height}\}$ . Both subgraphs are DAGs. Since  $A(G') \subseteq I$ , I is filling. But  $A(G'') \not\subseteq I$ , thus I' is not filling.

While in our simple examples it is easy to verify if a (sub-)graph is a DAG, in more involved settings the following algorithm could be used systematically.

A simple algorithm to check if G is a DAG uses the following algorithm.

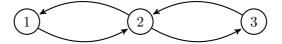
**Algorithm 2.15** (Directed path). Whenever one deals with directed paths in graphs, Dijkstra's Algorithm as described by himself in [3] is one of the most useful tools: it finds the shortest path (if there is any) between vertices. So let G be a directed graph on n vertices and let  $M_G$  be its adjacency matrix. There is a directed path from vertex i to vertex j if and only if  $M_G^k[i,j] > 0$  for some  $k \in \{1, \ldots, n-1\}$ .

**Algorithm 2.16** (Checking whether G is a DAG). Since in a DAG with n vertices, we have no path of length n, Algorithm 2.15 gives the following criterion: G is a DAG if and only if the sum of all traces of  $M_G^1, M_G^2, \ldots, M_G^{n-1}$  is 0. Other algorithms exist to identify if a directed graph is acyclic, see e.g. [4].

**Definition 2.17.** We say that a filling subset  $I \subseteq V$  is minimal (with respect to  $\subseteq$ ) if no proper subset of I is filling.

The following example shows that two different minimal filling subsets do not necessarily have the same number of elements:

**Example 2.18.** Let  $V = \{1, 2, 3\}$  and  $E = \{(1, 2), (2, 1)\} \cup \{(2, 3), (3, 2)\}.$ 



Then  $I = \{1,3\}$  is minimal filling, and also  $J = \{2\}$ , but I and J differ in cardinality. Note that, even if we could visit each vertex of the graph by starting at  $K = \{1\}$  and following the directed edges, K is not filling: The reason is that  $\{2\}$  is not completely determined by K alone (there is also edge pointing from  $\{3\}$  to  $\{2\}$  and  $\{3\}$  is only determined through  $\{2\}$ ). Put differently,  $K_1 = K \cup \text{Dtm}(\{1\}) = \{1\}$ , thus  $K_1 = K_2 = K_3 = \cdots = \{1\}$ . Therefore  $K = \{1\}$  is not filling.

Can we choose a minimal filling subsets such that it intersects every cycle (or every minimal cycle) at exactly 1 vertex? Unfortunately not:

**Example 2.19.** Let  $V = \{0, 1, 2\}$ , let

$$E = (V \times V) \setminus \{(k, k) : k \in \{0, 1, 2\}\},\$$

and let G = (V, E). Note that  $A(G) = \emptyset$ . If we take  $K = \{k\}$  for some  $k \in \{0, 1, 2\}$ , it is easily verified that  $\operatorname{Dtm}(K) = \emptyset$ . So  $K_n = K$  for all  $n \in \mathbb{N}$ , and therefore K is not filling. So if  $I \subseteq V$  is minimal filling, it contains at least 2 vertices of  $V = \{1, 2, 3\}$ , say  $\{k_1, k_2\} \subseteq I$  for  $k_1 \neq k_2 \in \{1, 2, 3\}$ . But then  $C = \{k_1, k_2\}$  is a minimal cycle by definition of E, and  $|I \cap C| > 1$ .

By Theorem 2.8 or Theorem 2.13 it is not hard to verify whether a given set of vertices is filling or not. However, in general there will be no simple algorithm to identify the smallest possible filling subset of V. The following greedy algorithm will, however, usually provide a good approximation.

## Algorithm 2.20.

(1) Identify the set I = A(G) of all vertices without incoming edges.

- (2) Choose all vertices  $W = \{w_1, \dots, w_m\}$  in all minimal cycles that do not intersect I. Denote by  $n_i$  the number of cycles intersected by  $w_i$ ,  $1 \le i \le m$ .
- (3) If  $W \neq \emptyset$ , pick any  $w_i \in W$  with maximal  $n_i$  and set  $I := I \cup \{w_i\}$
- (4) Go to Step 2 as long as W contains at least two elements. Otherwise stop with I as solution.

The algorithm will terminate after maximal n-1 iterations with n being the number of vertices. Due to its greedy nature, it might miss the optimal solution in certain hypothetical cases.

**Example 2.21.** To illustrate the algorithm and also the power of Theorems 2.8 and 2.13, take the graph G on the right hand side of Figure 1. There are no vertices without incoming edge, thus A(G) is empty. There are four minimal cycles that contain the following vertices.

- (1) {Age, Pregnant}
- (2) {Sex, Pregnant}
- (3) {Height, Age, Pregnant}
- (4) {Sex, Height, Age, Pregnant}

Consequently, after the first iteration of the algorithm up to Step 3, the candidate filling set I consists of the vertex "Pregnant" (hits highest number of cycles and A(G) is empty). Furthermore, W consists all vertices, so we start with a second iteration but now W is left empty and the algorithm stops. Thus, even if only the field "Pregnant" is entered, the remaining fields can be autofilled without further input. Of course there are also filling subsets without containing "Pregnant", e. g.  $I' = \{Age, Sex\}$  (apply Theorem 2.8). Figure 2 presents the two subgraphs without edges pointing to  $I = \{Pregnant\}$  (left picture) resp. to I' (right picture), illustrating the idea of Theorem 2.13.

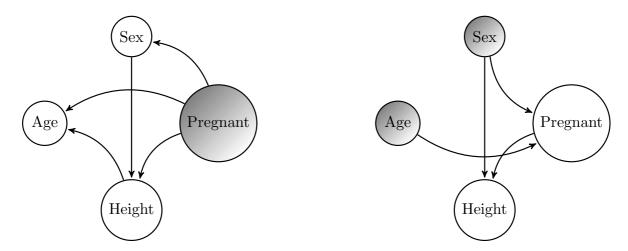


FIGURE 3. Subgraphs resulting from removing the incoming edges to  $I = \{\text{Pregnant}\}\ \text{resp.}\ I' = \{\text{Age, Sex}\}\ (\text{both highlighted})$  in the situation of Example 2.21. Since I resp. I' contain the vertices without incoming edges and since the subgraphs are acyclic, both I and I' are filling in the original graph without removed edges

# 3. FILLING WITH PARTIAL DETERMINATION

Given  $I \subseteq V$  and  $v \in V$  we say that I determines the vertex v with partial input if  $In(v) \cap I \neq \emptyset$ , that is if any argument in a replacement function is provided. For short, we say in that case that I p-determines v.

We denote the set of vertices p-determined by I by pDtm(I). Inductively we define the "p-determination closure" of  $I \subseteq V$  in the following way:

- $(1) \ I_0^{(p)} := I;$
- (1)  $I_0$  := I, (2) for  $n \ge 0$  set  $I_{n+1}^{(p)} := I_n^{(p)} \cup \mathrm{pDtm}(I_n)$ .

We say that  $I \subseteq V$  **p-fills** V if there is  $n \in \mathbb{N}$  such that  $I_n^{(p)} = V$ .

**Example 3.1.** In Example 2.18,  $I = \{1\}$  is p-filling (but not filling) because  $I_1 = \{1, 2\}$  and  $I_2 = \{1, 2, 3\}$ .

There is a first, trivial characterization of p-filling subsets of a graph G = (V, E):

**Proposition 3.2.** Let G = (V, E) be a self-loopless directed graph, and let  $I \subseteq V$ . Then the following are equivalent:

- I is p-filling;
- (2)  $A(G) \subseteq I$ , and for every vertex  $x \in V \setminus I$  there is a directed path from some  $j \in I$  to x.

However, we can do much better than that. Next, we specify in a mathematical way what it means to "collapse" points  $x, y \in V$  when there are directed paths between theses points in either direction.

**Definition 3.3.** If G = (V, E) is a directed graph and  $x, y \in V$  we say that x, y are strongly connected, in symbols  $x \simeq y$  if there exists a directed path from x to y and vice versa.

Again, it is straightforward to verify that  $\simeq$  is an equivalence relation on V. The set of elements equivalent to  $x \in V$  is denoted by  $[x]_{\simeq}$ , and we call it the *strongly connected component* (scc) containing x. The set of scc's on V with respect to  $\simeq$  is denoted by  $V/_{\simeq}$ . Note that by construction, every  $v \in V$  lies in a unique scc, the scc's are mutually disjoint, and all scc's are non-empty.

We put a directed graph structure on  $V/_{\sim}$  and set  $G/_{\sim}=(V/_{\sim},E/_{\sim})$  where

$$E/_{\simeq} = \{(C, D) \in V/_{\simeq} \times V/_{\simeq} : C \neq D \text{ and there are } c \in C, d \in D \text{ such that } (c, d) \in E\}.$$

The following is an elementary observation:

**Lemma 3.4.** If G is a directed graph,  $G/_{\simeq}$  is a DAG.

Before we show how strongly connected components come into play for finding p-filling sets, we need some basic observations.

# Lemma 3.5.

(1) If G = (V, E) is a DAG, then for all  $v \in V$  there is  $x \in A(G)$  and a directed path from x to v.

(2) Let G = (V, E) be any graph. If  $C, D \in G/_{\simeq}$  and there is a direct path in  $G/_{\simeq}$  from C to D then for all  $c \in C, d \in D$  there is a directed path in G from c to d.

Combining these observations lead us to the following:

**Proposition 3.6.** Let G = (V, E) be any graph, and let  $v \in V$ . Then there is a strongly connected component  $C_0 \in A(G/_{\simeq})$  such that for every  $c_0 \in C_0$  there is a directed path in G from  $c_0$  to v.

*Proof.* Let  $[v]_{\simeq}$  be the strongly connected component containing v. Lemma 3.4 says that  $G/_{\simeq}$  is a DAG. By Lemma 3.5 (1) there is  $C_0 \in A(G/_{\simeq})$  and a directed path in  $G/_{\simeq}$  from  $C_0$  to  $[v]_{\simeq}$ . Finally, Lemma 3.5 (2) implies that there is a directed path in G from any  $c_0 \in C_0$  to v.

**Theorem 3.7.** Let G = (V, E) be a directed graph,  $I \subseteq V$ , and consider the graph  $G/_{\simeq}$ . Then the following statements are equivalent:

- (1) I is p-filling;
- (2) I intersects every strongly connected component  $C \in A(G/_{\simeq})$ .

# Proof.

- (1)  $\Longrightarrow$  (2). Suppose  $C^* \in A(G/_{\simeq})$  and let  $I \subseteq V \setminus C^*$ . We show that I is not p-filling for G by establishing that  $I_n^{(p)} \cap C^* = \emptyset$  for all  $n \in \mathbb{N}$ . The statement is true for  $I_0^{(p)} = I$ . Assume that  $I_k^{(p)} \cap C^* = \emptyset$ . The fact that  $C^* \in A(G/_{\simeq})$  means by definition of  $A(\cdot)$  and by definition of  $E_{\simeq}$  that there is no edge in E coming into  $C^*$  from the outside or, more precisely, for all  $x \in V \setminus C^*$  and  $c \in C^*$  we have  $(x, c) \notin E$ . Therefore, for all  $c \in C^*$  we have  $c \notin I_{k+1}^{(p)}$ , that is  $C^* \cap I_{k+1}^{(p)} = \emptyset$ . This inductive argument proves that  $I_n^{(p)} \cap C^* = \emptyset$  for all  $n \in \mathbb{N}$ , so  $I_n^{(p)} \neq V$  for all n, and therefore I is not p-filling.
- (2)  $\Longrightarrow$  (1). Fix  $v \in V$ . By Proposition 3.2 we need to establish that if I intersects any strongly connected component (= element of  $A(G/_{\simeq})$ ), then there is a directed path from some  $i \in I$  to v. The proposition then implies that I is filling.

Use Proposition 3.6 to find  $C_0 \in A(G/_{\simeq})$  such that for every  $c_0 \in C_0$  there is a directed path in G from  $c_0$  to v. Since I intersects  $C_0$  by assumption, pick  $j_0 \in I \cap C_0$ . So there is a directed path in G from  $j_0 \in I$  to v. So by Proposition 3.2, I is filling because v was chosen arbitrarily.

**Algorithm 3.8.** In [7], Tarjan introduced an algorithm that identifies for any graph G = (V, E) its strongly connected components in time O(|V| + |E|).

Note that this theorem implies that minimal p-filling subsets intersect every member of  $A(G/_{\sim})$  at exactly one point. So this implies:

**Corollary 3.9.** All minimal (with respect to  $\subseteq$ ) p-filling subsets have the same cardinality.

This is in sharp contrast to the situation in the previous section, where minimal filling subsets can have different cardinalities (see Example 2.18).

Moreover, Theorem 3.7 gives an efficient algorithm to find minimally p-filling sets: Identify the strongly connected components that don't have an incoming edge (that is,  $A(G/_{\sim})$ ), and pick a vertex from each.

**Example 3.10.** In the introductory example of a weight-calculator with the vertices  $V = \{Age, Height, Sex\}$ , we could modify the replacement function

$$g(x,s) := \begin{cases} 162 + 16s & \text{if } x > 16, \\ \lfloor (x-1)/16 \cdot 130 + 30.5 \rfloor & \text{if } x \le 16 \end{cases}$$

for height z based on (non-missing) age x and sex s by a function that allows one of the two arguments to be missing, for instance by

$$g'(x,s) := \begin{cases} 162 + 16s & \text{if } (x > 16 \text{ or } x \text{ missing}) \text{ and } s \text{ non-missing,} \\ 170 & \text{if } x > 16 \text{ and } s \text{ missing,} \\ \lfloor (x-1)/16 \cdot 130 + 30.5 \rfloor & \text{if } x \leq 16. \end{cases}$$

Then, in our graph-theoretic autofill framework, the vertex "Height" would partially be determined by "Sex" and "Age" and "Age" (partially) determined by "Height". The graph on the left side of Figure 1 would be equivalent (under  $\simeq$ ) to the DAG  $G/_{\simeq}$  with the two strongly connected components {Sex} and {Age, Height} as vertices. By Theorem 3.7, any subset of V containing  $A(G/_{\simeq}) = \{\text{Sex}\}$  would be filling.

Generally, to use partial determination requires more effort to properly define the replacement functions compared to complete determination as these functions also need to treat missing input in a reasonable way. However, usually a smaller subset of input values is required to fill the remaining values.

# 4. Conclusion

Loopless directed graphs turn out to be the ideal framework for representing different kinds of "inference" or "implication". In this article, we used directed graphs in the context of missing values in vectors. In web forms used today, the set of mandatory fields is fixed. Our approach offers something new: flexible and smart filling of missing values.

### References

- [1] R. Diestel. **Graph Theory**, Springer Verlag, 2010.
- [2] D. Johnson. Finding all the elementary circuits of a directed graph, SIAM Journal on Computing 4(1): 77–84, 1975.
- [3] E. Dijkstra. A note on two problems in connexion with graphs, Numerische Mathematik 1: 269–271, 1959.
- [4] R. E. Tarajan. Edge-disjoint spanning trees and depth-first search, Acta Informatica 6(2): 171–18, 1976.
- [5] M. Saar-Tsechansky and F. Provost. Handling Missing Values when Applying Classification Models, Journal of Machine Learning Research 8: 1625–1657, 2007.
- [6] D. Schuurmans and R. Greiner. Learning to classify incomplete examples, Computational Learning Theory and Natural Learning Systems IV: Making Learning Systems Practical: 87–105, MIT Press, Cambridge MA, 1997.
- [7] R. E. Tarjan. Depth-first search and linear algorithms for graphs, SIAM Journal on Computing, 1(2): 146–160, 1972.

CONSULT AG, CH-8050 ZURICH, SWITZERLAND

 $E\text{-}mail\ address: \verb|michael.mayer@consultag.ch||$ 

Federal Office of Social Insurance, CH-3003 Bern, Switzerland

 $E\text{-}mail\ address: \verb|dominic.zypen@gmail.com||$