



Making, Updating, and Querying Causal Models using `CausalQueries`

Till Tietz 
WZB

Lily Medina
UC Berkeley

Georgiy Syunyaev 
Vanderbilt University

Macartan Humphreys 
WZB

Abstract

A guide to the R package `CausalQueries` for making, updating, and querying causal models

Keywords: causal models, stan, bayes.

1. Introduction: Causal models

`CausalQueries` is an R package that lets users make, update, and query causal models. Users provide a statement that reports a set of binary variables and the relations of causal ancestry between them: which variables are direct causes of other variables, given the other variables in the model. Once provided to `make_model()`, `CausalQueries` generates a parameter vector that fully describes a probability distribution over all possible types of causal relations between variables (“causal types”), given the causal structure. Given a prior over parameters and data over some or all nodes, `update_model()` deploys a `stan` model in order to generate a posterior distribution over causal models. The function `query_model()` can then be used to ask any causal query of the model, using either the prior distribution, the posterior distribution, or a user-specified candidate vector of parameters.

In the next section we provide a short motivating example. We then describe how the package relates to existing available software. Section 4 gives an overview of the statistical model behind the package. Sections 5, 6, and 7 then describe the main functionality for the major operations using the package. We provide further computation

details in the final section.

2. Motivating example

Before providing details on package functionality we illustrate these three core functions by showing how to use *CausalQueries* to replicate the analysis in [Chickering and Pearl \(1996\)](#) (see also [Humphreys and Jacobs \(2023\)](#)). [Chickering and Pearl \(1996\)](#) seek to draw inference on causal effects in the presence of imperfect compliance. We have access to an instrument Z (a randomly assigned prescription for cholesterol medication), which is a cause of X (treatment uptake) but otherwise unrelated to Y (cholesterol). We imagine we are interested in three specific queries. The first is the average causal effect of X on Y . The second is the average effect for units for which $X = 0$ and $Y = 0$. The last is the average treatment effect *for* “compliers”: units for which X responds positively to Z . Thus two of these queries are conditional queries, with one conditional on a counterfactual quantity.

Our data on Z , X , and Y is complete for all units and looks, in “compact form”, as follows:

```
R> data("lipids_data")
R>
R> lipids_data

#>   event strategy count
#> 1 Z0X0Y0      ZXY   158
#> 2 Z1X0Y0      ZXY    52
#> 3 Z0X1Y0      ZXY     0
#> 4 Z1X1Y0      ZXY    23
#> 5 Z0X0Y1      ZXY    14
#> 6 Z1X0Y1      ZXY    12
#> 7 Z0X1Y1      ZXY     0
#> 8 Z1X1Y1      ZXY    78
```

Note that in compact form we simply record the number of units (“count”) that display each possible pattern of outcomes on variables (“event”).¹

With *CausalQueries*, you can create the model, input data to update it, and then query the model for results.

```
R> make_model("Z -> X -> Y; X <-> Y") |>
+   update_model(lipids_data, refresh = 0) |>
+   query_model(query = "Y[X=1] - Y[X=0]",
+                 given = c("All", "X==0 & Y==0", "X[Z=1] > X[Z=0]"),
+                 using = "posteriors")
```

The output is a data frame with estimates, posterior standard deviations, and credibility intervals. For example the data frame produced by the code above is shown in [Table 1](#).

As we describe below the same basic procedure of making, updating, and querying models, can be used (up to computational constraints) for arbitrary causal models, for different types

¹The “strategy” column records the set of variables for which data has been recorded.

Table 1: Replication of Chickering and Pearl (1996).

query	given	mean	sd	cred.low.2.5%	cred.high.97.5%
$Y[X=1] - Y[X=0]$	-	0.56	0.10	0.38	0.73
$Y[X=1] - Y[X=0]$	$X==0 \ \& \ Y==0$	0.64	0.15	0.38	0.89
$Y[X=1] - Y[X=0]$	$X[Z=1] > X[Z=0]$	0.70	0.05	0.60	0.80

Rows 1 and 2 replicate results in Chickering and Pearl (1996); row 3 returns inferences for complier average effects.

of data structures, and for all causal queries that can be posed of the causal model.

3. Connections to existing packages

The literature on causal inference and its software ecosystem are rich and expansive; spanning the social and natural sciences as well as computer science and applied mathematics. In the interest of clarity we thus briefly contextualize **CausalQueries**' scope and functionality within the subset of the causal inference domain addressing the evaluation of causal queries on causal models encoded as directed acyclic graphs (DAGs) or structural equation models (SEMs). Table 2 provides an overview of relevant software and discusses key connections, advantages and disadvantages with respect to **CausalQueries**.

Table 2: Related software.

Software	Source	Language	Availability	Scope
causalnex	Beaumont, Horsburgh, Pilgerstorfer, Droth, Oentaryo, Ler, Nguyen, Ferreira, Patel, and Leong (2021)	Python	pip installable	<ul style="list-style-type: none"> causal structure learning querying marginal distributions discrete data
causaloptim	Michael C. Sachs and Gabriel (2023)	R	CRAN	<ul style="list-style-type: none"> bounding causal effects non-identified queries binary data
pclag	Kalisch, Mächler, Colombo, Maathuis, and Bühlmann (2012)	R	CRAN	<ul style="list-style-type: none"> causal structure learning ATEs under linear conditional expectations and no hidden selection

Software	Source	Language	Availability	Scope
autobounds	Duarte, Finkelstein, Knox, Mummolo, and Shpitser (2023)	Python	not readily installable	<ul style="list-style-type: none"> • bounding causal effects • partial identification • DAG canonicalization • binary data

One of the most comprehensive pieces of software in the causal modeling domain is **causalnex**. It provides a feature-rich and highly optimized set of tools to learn, update and query causal models using discrete data. While avoiding the rich model parameterization via principal strata (nodal types) employed by **CausalQueries**, allows **causalnex** to easily handle non-binary data and scale to large causal models; it substantially restricts the set of possible queries that can be evaluated and prior knowledge that can be specified on models. **causalnex** is in this sense akin to machine learning approaches to causal inference, focusing on causal structure learning in variable rich but potentially domain knowledge poor settings and the evaluation of simple queries over marginal distributions on learned DAGs. The rich model structure employed by **CausalQueries** by contrast makes it highly suited to problems where answers to complex causal queries are sought in relatively more domain knowledge abundant settings.

Like **causalnex**, **pclag** places particular emphasis on causal structure learning, utilizing the resultant DAGs to recover ATEs across all learned markov-equivalent classes implied by observed data that satisfy linearity of conditional expectations. This approach again is more restrictive than **CausalQueries** in the DAGs and particularly the queries it allows.

The software bearing the highest resemblance to **CausalQueries** with respect to model definition are **autobounds** and **causaloptim**. Dealing with binary causal models, their definitions of principal strata (nodal types) and the resultant set of causal relations on the DAG (causal types) are very close to those of **CausalQueries**. Differences in model definition arise with respect to disturbance terms and confounding being defined implicitly via main nodes and edges in **CausalQueries** vs explicitly via separate disturbance nodes in **autobounds** and **causaloptim**. While **CausalQueries** assumes canonical form for input DAGs, **autobounds** and **causaloptim** facilitate canonicalization. The essential difference between the methods; however, lies in their approach to evaluating queries.

Both **autobounds** and **causaloptim** build on seminal approaches in [Balke and Pearl \(1997\)](#) to construct bounds of queries, using constrained polynomial and linear optimization respectively. In contrast, **CausalQueries** utilizes Bayesian inference to generate a posterior over the causal model which is then queried (consistent with [Chickering and Pearl \(1996\)](#) and [Zhang, Tian, and Bareinboim \(2022\)](#)). A key difference then is the target of inference. The polynomial and linear programming approach to querying is in principle suited to handling larger causal models, though given their similarity in model parameterization, **autobounds**, **causaloptim** and **CausalQueries** face similar constraints induced by parameter spaces expanding rapidly with model size. The Bayesian approach to model updating and querying holds the efficiency advantage that a model can be updated once and queried infinitely, while expensive optimization runs are required for each separate query in **autobounds** and

`causaloctim`.

Summarizing, the particular strength of **CausalQueries** is to allow users to specify arbitrary DAGs, arbitrary queries over nodes in those DAGs, and use the same canonical procedure to form Bayesian posteriors over those queries whether or not the queries are identified. Thus in principle if researchers are interested in learning about a quantity like the local average treatment effect and their model in fact satisfies the conditions in [Angrist, Imbens, and Rubin \(1996\)](#), then updating will recover valid estimates even if researchers are unaware that the local average treatment effect is identified and are ignorant of the estimation procedure proposed by [Angrist et al. \(1996\)](#).

There are two broad limitations on the sets of models handled natively by **CausalQueries**. First **CausalQueries** is designed for models with a relatively small number of binary nodes. Because there is no compromise made on the space of possible causal relations implied by a given model, the parameter space grows very rapidly with the complexity of the causal model. The complexity also depends on the causal structure and grows rapidly with the number of parents affecting a given child. A chain model of the form $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E$ has just 40 parameters. A model in which A, B, C, D are all direct ancestors of E has 65544 parameters. Moving from binary to non binary nodes has similar effects. The restriction to binary nodes is for computational and not conceptual reasons. In fact it is possible to employ **CausalQueries** to answer queries from models with non binary nodes but in general the computational costs make analysis of these models prohibitive.²

Second, the package is geared towards learning about populations from samples of units that are independent of each other and are independently randomly sampled from populations. Thus the basic set up does not address problems of sampling, clustering, hierarchical structures, or purposive sampling. The broader framework can however be used for these purposes (see section 9.4 of [Humphreys and Jacobs \(2023\)](#)). The targets of inference are usually case level quantities or population quantities and **CausalQueries** is not well suited for estimating sample quantities.

4. Statistical model

The core conceptual framework is described in Pearl’s *Causality* ([Pearl 2009](#)) but can be summarized as follows (using the notation used in [Humphreys and Jacobs \(2023\)](#)):

Definition 1 A “*causal model*” is:

1. an ordered collection of “endogenous nodes” $Y = \{Y_1, Y_2, \dots, Y_n\}$
2. an ordered collection of “exogenous nodes” $\Theta = \{\theta^{Y_1}, \theta^{Y_2}, \dots, \theta^{Y_n}\}$
3. a collection of functions $F = \{f_{Y_1}, f_{Y_2}, \dots, f_{Y_n}\}$ specifying, for each j , how outcome y_j depends on θ_j and realizations of endogenous nodes prior to j .
4. a probability distribution over Θ , λ .

²For more on computation constraints and strategies to update and query large models see the associated package **CausalQueriesTools**. The core approach used here is to divide large causal models into modules, update on modules and reassemble to pose queries.

In the usual case *CausalQueries* takes endogenous nodes to be binary.³ When we specify a causal structure we specify which endogenous nodes are (possibly) direct causes of a node, Y_j , given other nodes in the model. These nodes are called the parents of Y_j , PA_j (we use upper case PA_j to indicate the collection of nodes and lower case pa_j to indicate a particular set of values that these nodes might take on). With discrete valued nodes, it is possible to identify all possible ways that a node might respond to its parents. We refer to the ways that a node responds as “nodal type.” The set of nodal types corresponds to principal strata familiar, for instance, in the study of instrumental variables (Frangakis and Rubin 2002).

If node Y_i can take on k_i possible values then the set of possible values that can be taken on by parents of j is $m := \prod_{i \in PA_j} k_i$. Then there are k_j^m different ways that a node might respond to its parents. In the case of binary nodes this becomes $2^{(2^{|PA_j|})}$. Thus for an endogenous node with no parents there are 2 nodal types, for a binary node with one binary parent there are four types, for a binary node with 2 parents there are 16, and so on.

The set of all possible causal reactions of a given unit to all possible values of parents is then given by its collection of nodal types at each node. We call this collection a unit’s “causal type”, θ .

The approach used by *CausalQueries* is to let the domain of θ^{Y_j} be coextensive with the number of nodal types for Y_j . Function f^j then determines the value of y by simply reporting the value of Y_j implied by the nodal type and the values of the parents of Y_j . Thus if $\theta_{pa_j}^j$ is the value for j when parents have values pa_j , then we have simply that $f_{y_j}(\theta^j, pa_j) = \theta_{pa_j}^j$. The practical implication is that, given the causal structure, learning about the model reduces to learning about the distribution, λ , over the nodal types.

In cases in which there is no unobserved confounding, we take the probability distributions over the nodal types for different nodes to be independent: $\theta^i \perp\!\!\!\perp \theta^j, i \neq j$. In this case we use a categorical distribution to specify the $\lambda^{j'} := \Pr(\theta^j = \theta^{j'})$. From independence then we have that the probability of a given causal type θ' is simply $\prod_{i=1}^n \lambda^{i'}$.

In cases in which there is confounding, the logic is essentially the same except that we need to specify enough parameters to capture the joint distribution over nodal types for different nodes.

We make use of the causal structure to simplify. As an example, for the lipids model, the full joint distribution of nodal types can be simplified as in Equation 1.

$$\Pr(\theta^Z = \theta_1^Z, \theta^X = \theta_{10}^X, \theta^Y = \theta_{11}^Y) = \Pr(\theta^Z = \theta_1^Z) \Pr(\theta^X = \theta_{10}^X) \Pr(\theta^Y = \theta_{11}^Y | \theta^X = \theta_{10}^X) \quad (1)$$

And so for this model λ would include parameters that represent $\Pr(\theta^Z)$ and $\Pr(\theta^X)$ but also the conditional probability $\Pr(\theta^Y | \theta^X)$:

$$\Pr(\theta^Z = \theta_1^Z, \theta^X = \theta_{10}^X, \theta^Y = \theta_{11}^Y) = \lambda_1^Z \lambda_{10}^X \lambda_{11}^{Y|\theta_{10}^X} \quad (2)$$

³*CausalQueries* can be used also to analyse non binary data though with a cost of greatly increased complexity. See section 9.4.1 of Humphreys and Jacobs (2023) for an approach that codes non binary data as a profile of outcomes on multiple binary nodes.

Representing beliefs *over causal models* thus requires specifying a probability distribution over λ . This might be a degenerate distribution if users want to specify a particular model. **CausalQueries** allows users to specify parameters, α of a Dirichlet distribution over λ . If all entries of α are 0.5 this corresponds to Jeffrey’s priors. The default behavior is for **CausalQueries** to assume a uniform distribution – that is, that all nodal types are equally likely – which corresponds to α being a vector of 1s.

Updating is then done with respect to beliefs over λ . In the Bayesian approach we have simply:

$$p(\lambda'|D) = \frac{p(D|\lambda')p(\lambda')}{\int_{\lambda''} p(D|\lambda'')p(\lambda'')}$$

$p(D|\lambda')$ is calculated under the assumption that units are exchangeable and independently drawn. In practice this means that the probability that two units have causal types θ_i and θ_j is simply $\lambda'_i\lambda'_j$. Since a causal type fully determines an outcome vector $d = \{y_1, y_2, \dots, y_n\}$, the probability of a given outcome (“event”), w_d , is given simply by the probability that the causal type is among those that yield outcome d . Thus from λ' we can calculate a vector of event probabilities, $w(\lambda)$, for each vector of outcomes, and under independence we have:

$$D \sim \text{Multinomial}(w(\lambda), N)$$

Thus for instance in the case of a $X \rightarrow Y$ model, and letting w_{xy} denote the probability of a data type $X = x, Y = y$, the event probabilities are:

$$w(\lambda) = \begin{cases} w_{00} &= \lambda_0^X(\lambda_{00}^Y + \lambda_{01}^Y) \\ w_{01} &= \lambda_0^X(\lambda_{11}^Y + \lambda_{10}^Y) \\ w_{10} &= \lambda_1^X(\lambda_{00}^Y + \lambda_{10}^Y) \\ w_{11} &= \lambda_1^X(\lambda_{11}^Y + \lambda_{01}^Y) \end{cases}$$

For concreteness: Table 3 illustrates key values for the lipids model. We see here that we have two types for node Z , four for X (representing the strata familiar from instrumental variables analysis: never takers, always takers, defiers, and compliers) and 4 for Y . For Z and X we have parameters corresponding to probability of these nodal types. For instance $Z.0$ is the probability that $Z = 1$. $Z.1$ is the complementary probability that $Z = 1$. Things are little more complicated for distributions on nodal types for Y however: because of confounding between X and Y we have parameters that capture the conditional probability of the nodal types for Y *given* the nodal types for X . We see there are four sets of these parameters.

The final column shows a sample set of parameter values. Together the parameters describe a full joint probability distribution over types for Z , X and Y that is faithful to the graph.

From these we can calculate the probability of each data type. For instance the probability of data type $Z = 0, X = 0, Y = 0$ is:

$$w_{000} = \Pr(Z = 0, X = 0, Y = 0) = \lambda_0^Z \left(\lambda_{00}^X(\lambda_{00}^{Y|\lambda_{00}^X} + \lambda_{01}^{Y|\lambda_{00}^X}) + \lambda_{01}^X(\lambda_{00}^{Y|\lambda_{01}^X} + \lambda_{01}^{Y|\lambda_{01}^X}) \right)$$

Table 3: Nodal types and parameters for Lipids model

node	nodal_type	param_set	param_names	param_value
Z	0	Z	Z.0	0.71
Z	1	Z	Z.1	0.29
X	00	X	X.00	0.36
X	10	X	X.10	0.03
X	01	X	X.01	0.51
X	11	X	X.11	0.10
Y	00	Y.X.00	Y.00_X.00	0.43
Y	10	Y.X.00	Y.10_X.00	0.08
Y	01	Y.X.00	Y.01_X.00	0.34
Y	11	Y.X.00	Y.11_X.00	0.15
Y	00	Y.X.01	Y.00_X.01	0.43
Y	10	Y.X.01	Y.10_X.01	0.05
Y	01	Y.X.01	Y.01_X.01	0.39
Y	11	Y.X.01	Y.11_X.01	0.13
Y	00	Y.X.10	Y.00_X.10	0.24
Y	10	Y.X.10	Y.10_X.10	0.45
Y	01	Y.X.10	Y.01_X.10	0.12
Y	11	Y.X.10	Y.11_X.10	0.19
Y	00	Y.X.11	Y.00_X.11	0.61
Y	10	Y.X.11	Y.10_X.11	0.11
Y	01	Y.X.11	Y.01_X.11	0.03
Y	11	Y.X.11	Y.11_X.11	0.25

In practice **CausalQueries** uses a matrix `parmap` that maps from parameters into data types.

The value of the **CausalQueries** package is to allow users to specify *arbitrary* models of this form, figure out all the implied nodal types and causal types, and then update given priors and data by calculating event probabilities implied by all possible parameter vectors and in turn the likelihood of the data given the model. In addition, the package allows for arbitrary querying of a model to assess the values of estimands of interest that a re function of the values or counterfactual values of nodes conditional on values or counterfactual values of nodes.

In the next sections we review key functionality for making, updating and querying causal models.

5. Making models

A model is defined in one step in **CausalQueries** using a `dagitty` syntax in which the structure of the model is provided as a statement.

For instance:

```
R> model <- make_model("X -> M -> Y <- X")
```

The statement in quotes, "X -> M -> Y <- X", provides the names of nodes. An arrow ("->" or "<-") connecting nodes indicates that one node is a potential cause of another, i.e. whether a given node is a "parent" or "child" of another.

Formally a statement like this is interpreted as:

1. Functional equations:

- $Y = f(M, X, \theta^Y)$
- $M = f(X, \theta^M)$
- $X = \theta^X$

2. Distributions on Θ :

- $\Pr(\theta^i = \theta_k^i) = \lambda_k^i$

3. Independence assumptions:

- $\theta_i \perp\!\!\!\perp \theta_j, i \neq j$

where function f maps from the set of possible values of the parents of i to values of node i given θ^i as described above.

In addition, as we did in the [Chickering and Pearl \(1996\)](#) example, it is possible to use two headed arrows (<->) to indicate "unobserved confounding", that is, the presence of an unobserved variable that might influence observed variables. In this case condition 3 above is relaxed and the exogenous nodes associated with confounded variables have a joint distribution. We describe how this is done in greater detail in section [Section 5.3.2](#).

5.1. Graphing

Plotting the model can be useful to check that you have defined the structure of the model correctly. *CausalQueries* provides simple graphing tools that draw on functionality from the *dagitty*, *ggplot2*, and *ggdag* packages.

Once defined, a model can be graphed by calling the `plot()` method defined for the objects with class `causal_model` produced by `make_model()` function.

```
R> make_model("X -> M -> Y <- X; Z -> Y") |>
+ plot()
```

Alternatively you can provide a number of options to the `plot()` call that will be passed to `CausalQueries::plot_dag` via the method.

```
R> make_model("X -> M -> Y <- X; Z -> Y") |>
+ plot(x_coord = 1:4,
+       y_coord = c(1.5,2,1,2),
+       textcol = "white",
+       textsize = 3,
+       shape = 18,
+       nodecol = "grey",
+       nodesize = 12)
```

The graphs produced by the two calls above are shown in Figure 1. In both cases the resulting plot will have class `c("gg", "ggplot")` and so will accept any additional modifications available via the *ggplot2* package.

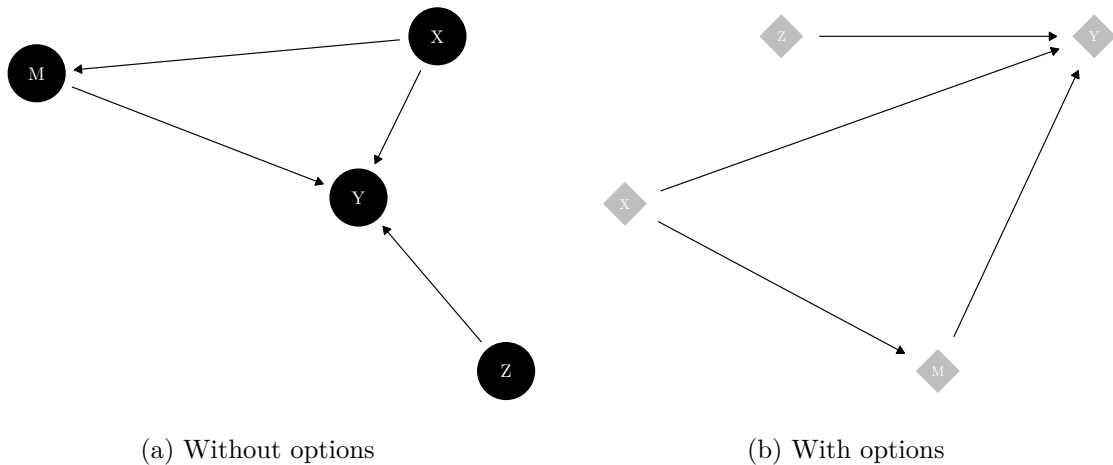


Figure 1: Examples of model graphs.

5.2. Model characterization

When a model is defined, a set of objects is generated. These are the key quantities that are used for all inference. The table below summarizes the core components of a model, providing a brief explanation for each one.

The first element is a **statement** which defines how the nodes in the model are related,

specified by the user using **dagitty** syntax. The second element, **dag**, is a dataframe that outlines the parent-child relationships within the model. The element **nodes** is simply a list of the names of the nodes in the model. Lastly, **parents_df**, is a table listing the nodes, indicating if they are “root” nodes (nodes with no parents among the set of specified nodes), and showing how many parents each node has.

The model includes additional elements, **nodal_types**, **parameters_df**, and **causal_types**, which we explain in detail later.

Table 4: Core Elements of a Causal Model.

Element	Description
statement	A character string that describes directed causal relations between variables in a causal model, where arrows denote that one node is a potential cause of another.
dag	A data frame with columns ‘parent’ and ‘children’ indicating how nodes relate to each other.
nodes	A list containing the nodes in the model.
parents_df	A table listing nodes, whether they are root nodes or not, and the number of parents they have.
nodal_types	A list with the nodal types in the model. See Section 5.2.2 for more details.
parameters_df	A data frame linking the model’s parameters with the nodal types of the model, as well as the family to which they belong. See Section 5.2.1 for more details.
causal_types	A data frame listing causal types and the nodal types that produce them. (See Causal Types Section)

After updating a model, two additional components are attached to it:

- A posterior distribution of the parameters in the model, generated by **stan**. This distribution reflects the updated parameter values.
- A list of other optional objects, **stan_objects**. The **stan_objects** can include the **stanfit** object and distributions over nodal types and event probabilities (**w**).

The table below summarizes the objects attached to the model after updating.

Table 5: Additional Elements.

Element	Description
posterior_distribution	The posterior distribution of the updated parameters generated by stan .
stan_objects	A list of additional objects (see next rows).
data	The data used for updating the model, always included in stan_objects .

Table 6: Example of Parameters Data Frame

param_names	node	gen	param_set	nodal_type	given	param_value	priors
X.0	X	1	X	0		0.50	1
X.1	X	1	X	1		0.50	1
Y.00	Y	2	Y	00		0.25	1
Y.10	Y	2	Y	10		0.25	1
Y.01	Y	2	Y	01		0.25	1
Y.11	Y	2	Y	11		0.25	1

Element	Description
type_distribution	The updated distribution of the nodal types, appended to stan_objects by default.
w	A mapping from parameters to event probabilities, optionally appended to stan_objects .
stan_fit	The stanfit object generated by stan, optionally appended to stan_objects .

Parameters data frame

When a model is created, *CausalQueries* attaches a “parameters data frame” which keeps track of model parameters, which belong together in a family, and how they relate to causal types. This becomes especially important for more complex models with confounding that might involve more complicated mappings between parameters and nodal types. In the case with no confounding the nodal types *are* the parameters; in cases with confounding there are generally more parameters than nodal types.

For instance:

```
R> make_model("X -> Y")$parameters_df
```

prints the parameters data frame shown in Table 6. Each row in the data frame corresponds to a single parameter.

The columns of the parameters data frame are understood as follows:

- **param_names** gives the name of the parameter, in shorthand. For instance the parameter $\lambda_0^X = \Pr(\theta^X = \theta_0^X)$ has **par_name** X.0.
- **param_value** gives the (possibly default) parameter values (probabilities).
- **param_set** indicates which parameters group together to form a simplex. The parameters in a set have parameter values that sum to 1. In this example $\lambda_0^X + \lambda_1^X = 1$.
- **node** indicates the node associated with the parameter.
- **nodal_type** indicates the nodal types associated with the parameter.
- **gen** indicates the place in the partial causal ordering (generation) of the node associated with the parameter

- `priors` gives (possibly default) Dirichlet priors arguments for parameters in a set. Values of 1 (.5) for all parameters in a set implies uniform (Jeffrey's) priors over this set.

Below we will see examples where the parameters data frame helps keep track of parameters that are created when confounding is added to a model.

Nodal types

As described above, two units have the same *nodal type* at node Y , θ^Y , if their outcome at Y responds in the same ways to parents of Y .

A binary node with k binary parents has 2^{2^k} nodal types. The reason is that with k parents, there are 2^k possible values of the parents and so 2^{2^k} ways to respond to these possible parental values. As a convention we say that a node with no parents has two nodal types (0 or 1).

When a model is created the full set of nodal types is identified. These are stored in the model. The labels for these nodal types indicate how the unit responds to values of parents.

For instance, consider the model with two parents $X \rightarrow Y \leftarrow M$. In such a case, the nodal types of Y will have subscripts with four digits, with each digit representing one of the possible combinations of values that Y can take, given the values of its parents X and M . These combinations include the value of Y when:

- $X = 0$ and $M = 0$,
- $X = 0$ and $M = 1$,
- $X = 1$ and $M = 0$,
- $X = 1$ and $M = 1$.

As the number of parents increases, keeping track of what each digit represents becomes more difficult. For instance, if Y had three parents, its nodal types would have subscripts of eight digits, each associated with the value that Y would take for each combination of the three parents. The `interpret_type` function provides a clear map to identify what each digit in the subscript represents. See the examples below for models with two and three parents.

Interpretations are automatically provided as part of the model object. A user can see them like this.

```
R> make_model("X -> Y")$nodal_types

#> $X
#> [1] "0" "1"
#>
#> $Y
#> [1] "00" "10" "01" "11"
#>
#> attr("interpret")
#> attr("interpret")$X
#>   node position display interpretation
#> 1     X         NA      X0           X = 0
#> 2     X         NA      X1           X = 1
```

```
#>
#> attr("interpret")$Y
#>   node position display interpretation
#> 1     Y          1   Y[*]*          Y | X = 0
#> 2     Y          2   Y*[*]          Y | X = 1
```

The `interpret_type` function can also be called by the user to obtain interpretations for the nodal types of each node in the model.

```
R> interpretations <-
+   make_model("X -> Y <- M; W -> Y") |>
+   interpret_type()
R>
R> interpretations$Y

#>   node position      display      interpretation
#> 1     Y          1 Y[*]***** Y | M = 0 & W = 0 & X = 0
#> 2     Y          2 Y*[*]***** Y | M = 1 & W = 0 & X = 0
#> 3     Y          3 Y**[*]***** Y | M = 0 & W = 1 & X = 0
#> 4     Y          4 Y***[*]***** Y | M = 1 & W = 1 & X = 0
#> 5     Y          5 Y****[*]*** Y | M = 0 & W = 0 & X = 1
#> 6     Y          6 Y*****[*]** Y | M = 1 & W = 0 & X = 1
#> 7     Y          7 Y*****[*]* Y | M = 0 & W = 1 & X = 1
#> 8     Y          8 Y*****[*] Y | M = 1 & W = 1 & X = 1
```

Causal types

Causal types are collections of nodal types. Two units are of the same *causal type* if they have the same nodal type at every node. For example in a $X \rightarrow M \rightarrow Y$ model, $\theta = (\theta_0^X, \theta_{01}^M, \theta_{10}^Y)$ is a type that has $X = 0$, M responds positively to X , and Y responds positively to M .

When a model is created, the full set of causal types is identified. These are stored in the model object:

```
R> lipids_model$causal_types |> head()

#>           Z   X   Y
#> Z0.X00.Y00 0 00 00
#> Z1.X00.Y00 1 00 00
#> Z0.X10.Y00 0 10 00
#> Z1.X10.Y00 1 10 00
#> Z0.X01.Y00 0 01 00
#> Z1.X01.Y00 1 01 00
```

In the lipids model there are $2 \times 4 \times 4 = 32$ causal types. A model with n_j nodal types at node j has $\prod_j n_j$ causal types. Thus the set of causal types can be large.

Knowledge of a causal type tells us what values a unit would take, on all nodes, whether or not there are interventions. For example for a model $X \rightarrow M \rightarrow Y$ a type $\theta = (\theta_0^X, \theta_{01}^M, \theta_{10}^Y)$ would imply data $(X = 0, M = 0, Y = 1)$ absent any intervention. (The converse of this, of course,

is the key to updating: observation of data ($X = 0, M = 0, Y = 1$) result in more weight placed on θ_0^X , θ_{01}^M , and θ_{10}^Y .) The general approach used by **CausalQueries** for calculating outcomes from causal types is given in section `#sec-propagation`.

Parameter matrix

The parameters data frame keeps track of parameter values and priors for parameters but it does not provide a mapping between parameters and the probability of causal types.

The parameter matrix (P matrix) can be added to the model to provide this mapping. The P matrix has a row for each parameter and a column for each causal type. For instance:

```
R> make_model("X -> Y") |> get_parameter_matrix()

#>
#> Rows are parameters, grouped in parameter sets
#>
#> Columns are causal types
#>
#> Cell entries indicate whether a parameter probability is used
#> in the calculation of causal type probability
#>
#>      X0.Y00 X1.Y00 X0.Y10 X1.Y10 X0.Y01 X1.Y01 X0.Y11 X1.Y11
#> X.0       1      0      1      0      1      0      1      0
#> X.1       0      1      0      1      0      1      0      1
#> Y.00      1      1      0      0      0      0      0      0
#> Y.10      0      0      1      1      0      0      0      0
#> Y.01      0      0      0      0      1      1      0      0
#> Y.11      0      0      0      0      0      0      1      1
#>
#>
#> param_set  (P)
#>
```

The probability of a causal type is given by the product of the parameter values for parameters whose row in the P matrix contains a 1.

Later (e.g. `?@tbl-confound-params-mat`) we will see examples where the P matrix helps keep track of parameters that are created when confounding is added to a model.

The parameter matrix is generated on the fly as needed, but it can also be added to the model using `set_parameter_matrix()`, which can sometimes be useful to speed up operations:

```
R> model <- model |> set_parameter_matrix()
```

5.3. Tailoring models

Setting restrictions

When a model is defined, the complete set of possible causal relations is worked out. This set can be very large.

Sometimes for theoretical or practical reasons it is useful to constrain the set of types. In *CausalQueries* this is done at the level of nodal types, with restrictions on causal types following from restrictions on nodal types.

To illustrate, in analyses of data with imperfect compliance, like we saw in our motivating *lipids* model example, it is common to impose a monotonicity assumption: that X does not respond negatively to Z . This is one of the conditions needed to interpret instrumental variables estimates as (consistent) estimates of the complier average treatment effect.

In *CausalQueries* we impose this assumption thus

```
R> model_restricted <-
+ make_model("Z -> X -> Y; X <-> Y") |>
+ set_restrictions("X[Z=1] < X[Z=0]")
```

in words: we restrict by removing types for which X is decreasing in Z .

We can use `get_parameter_matrix(model)` to view the resulting parameter matrix in which both the set of parameters and the set of causal types are restricted.

```
R> get_parameter_matrix(model)
```

Here and in general, setting restrictions typically involves using causal syntax; see Section [7.2](#) for a guide the syntax used by *CausalQueries*.

Note:

- Restrictions have to operate on nodal types: restrictions on *levels* of endogenous nodes are not allowed. This, for example, will fail: `make_model("X -> Y") |> set_restrictions(statement = "(Y == 1)")`. The reason is that it requests a correlated restriction on nodal types for X and Y which involves undeclared confounding.
- Restrictions implicitly assume fixed values for *all* parents of a node. For instance: `make_model("A -> B <- C") |> set_restrictions("B[C=1]==1")` is interpreted as shorthand for the restriction " $B[C = 1, A = 0] == 1 \mid B[C = 1, A = 1] == 1$ ".
- To place restrictions on multiple nodes at the same time, provide these as a vector of restrictions. This is not permitted: `set_restrictions("Y[X=1]==1 & X==1")`, since it requests correlated restrictions. This however is allowed: `set_restrictions(c("Y[X=1]==1", "X==1"))`.
- Use the `keep` argument to indicate whether nodal types should be dropped (default) or retained.
- Restrictions based on causal statements can use helpers, `decreasing()`, `increasing()`, `'complements()'` and similar. For instance `set_restrictions(decreasing("X", "Y"))`
- Restrictions can be set using nodal type labels.

```
R> make_model("S -> C -> Y <- R <- X; X -> C -> R") |>
+ set_restrictions(labels = list(C = "1000", R = "0001", Y = "0001"),
```



```
+ keep = TRUE)
```

- Wild cards can be used in nodal type labels:

```
R> make_model("X -> Y") |>
+ set_restrictions(labels = list(Y = "?0"))
```

- in models with confounding restrictions can be added to nodal types conditional on the values of other nodal types; this is done using a “given” argument.

```
R> make_model("X -> Y; X <-> Y") |>
+ set_restrictions(labels = list(X = '0', Y = '00'),
+ given = c(NA, 'X.1')) |>
+ get_parameter_matrix()
```

	X1.Y00	X1.Y10	X1.Y01	X1.Y11
X.1	1	1	1	1
Y.00_X.1	1	0	0	0
Y.10_X.1	0	1	0	0
Y.01_X.1	0	0	1	0
Y.11_X.1	0	0	0	1

Allowing confounding

(Unobserved) confounding between two nodes arises when the nodal types for the nodes are not independently distributed.

In the $X \rightarrow Y$ graph, for instance, there are 2 nodal types for X and 4 for Y . There are thus 8 joint nodal types (or causal types):

Table 7: Nodal Types in $X \rightarrow Y$ Model.

$\theta^Y \setminus \theta^X$	0	1	Σ
00	$\Pr(\theta_0^X, \theta_{00}^Y)$	$\Pr(\theta_1^X, \theta_{00}^Y)$	$\Pr(\theta_{00}^Y)$
10	$\Pr(\theta_0^X, \theta_{10}^Y)$	$\Pr(\theta_1^X, \theta_{10}^Y)$	$\Pr(\theta_{10}^Y)$
01	$\Pr(\theta_0^X, \theta_{01}^Y)$	$\Pr(\theta_1^X, \theta_{01}^Y)$	$\Pr(\theta_{01}^Y)$
11	$\Pr(\theta_0^X, \theta_{11}^Y)$	$\Pr(\theta_1^X, \theta_{11}^Y)$	$\Pr(\theta_{11}^Y)$
Σ	$\Pr(\theta_0^X)$	$\Pr(\theta_1^X)$	1

This table has eight interior elements and so an unconstrained joint distribution would have 7 degrees of freedom. A no confounding assumption means that $\Pr(\theta^X | \theta^Y) = \Pr(\theta^X)$, or $\Pr(\theta^X, \theta^Y) = \Pr(\theta^X) \Pr(\theta^Y)$. In this case we just put a distribution on the marginals and there would be 3 degrees of freedom for Y and 1 for X , totaling 4 rather than 7.

```
R> confounded <- make_model("X -> Y ; X <-> Y")
R>
R> confounded$parameters_df
```

Table 8: Parameters Data Frame for Model with Confounding.

param_names	node	gen	param_set	nodal_type	given	param_value	priors
X.0	X	1	X	0		0.50	1
X.1	X	1	X	1		0.50	1
Y.00_X.0	Y	2	Y.X.0	00	X.0	0.25	1
Y.10_X.0	Y	2	Y.X.0	10	X.0	0.25	1
Y.01_X.0	Y	2	Y.X.0	01	X.0	0.25	1
Y.11_X.0	Y	2	Y.X.0	11	X.0	0.25	1
Y.00_X.1	Y	2	Y.X.1	00	X.1	0.25	1
Y.10_X.1	Y	2	Y.X.1	10	X.1	0.25	1
Y.01_X.1	Y	2	Y.X.1	01	X.1	0.25	1
Y.11_X.1	Y	2	Y.X.1	11	X.1	0.25	1

Table 9: Parameter Matrix for Model with Confounding.

	X0.Y00	X1.Y00	X0.Y10	X1.Y10	X0.Y01	X1.Y01	X0.Y11	X1.Y11
X.0	1	0	1	0	1	0	1	0
X.1	0	1	0	1	0	1	0	1
Y.00_X.0	1	0	0	0	0	0	0	0
Y.10_X.0	0	0	1	0	0	0	0	0
Y.01_X.0	0	0	0	0	1	0	0	0
Y.11_X.0	0	0	0	0	0	0	1	0
Y.00_X.1	0	1	0	0	0	0	0	0
Y.10_X.1	0	0	0	1	0	0	0	0
Y.01_X.1	0	0	0	0	0	1	0	0
Y.11_X.1	0	0	0	0	0	0	0	1

Table 8 shows the parameters data frame produced by the code above. We see here that there are now two parameter families for parameters associated with the node Y . Each family captures the conditional distribution of Y 's nodal types, given X . For instance the parameter $Y01_X.1$ can be interpreted as $\Pr(\theta^Y = \theta_{01}^Y | X = 1)$.

To see exactly how the parameters map to causal types we can view the parameter matrix:

```
R> get_parameter_matrix(confounded)
```

The output is shown in Table 9. Importantly, the P matrix works as before, despite confounding. We can assess the probability of causal types by multiplying the probabilities of the constituent parameters.

Setting Priors

Priors on model parameters can be added to the parameters data frame. The priors are interpreted as “alpha” arguments for a Dirichlet distribution. The Dirichlet distribution is a probability distribution over an $n - 1$ dimensional unit simplex. It can be thought of as

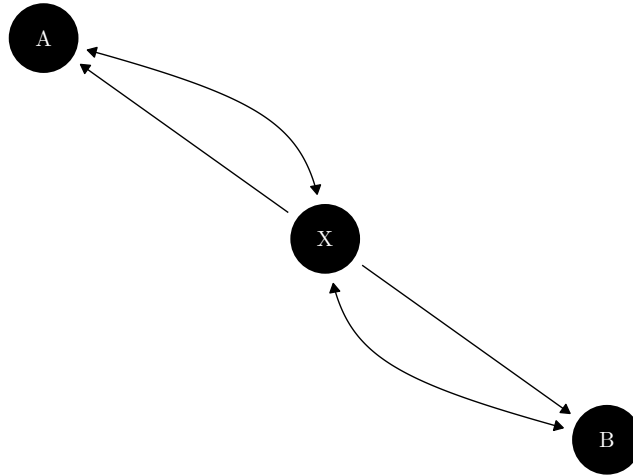


Figure 2: Graph of Model with Confounding.

a generalization of the Beta distribution and is parameterized by an n -dimensional positive vector α . Thus for example a Dirichlet with $\alpha = (1, 1, 1, 1, 1)$ gives a probability distribution over all non negative 5-dimensional vectors that sum to 1, e.g. $(0.1, 0.1, 0.1, 0.1, 0.6)$ or $(0.1, 0.2, 0.3, 0.3, 0.1)$. This particular value for α implies that all such vectors are equally likely. Other values for α can be used to control the expectation for each dimension as well as certainty. Thus for instance the vector $\alpha = (100, 1, 1, 1, 100)$ would result in more weight on distributions that are close to $(.5, 0, 0, 0, .5)$.

In **CausalQueries** priors are generally specified over the distribution of nodal types (or over the conditional distribution of nodal types, when there is confounding). Thus for instance in an $X \rightarrow Y$ model we have one Dirichlet distribution over the two types for θ^X and one Dirichlet distribution over the four types for θ^Y .

To retrieve the model's priors we can run the following code:

```
R> make_model("X -> Y") |> get_priors()

#>   X.0   X.1 Y.00 Y.10 Y.01 Y.11
#>    1     1     1     1     1     1
```

Here the priors have not been specified and so they default to 1, which corresponds to uniform priors. Alternatively you could set Jeffreys priors using **set_priors** as follows:

```
R> make_model("X -> Y") |>
+   set_priors(distribution = "jeffreys") |>
+   get_priors()

#>   X.0   X.1 Y.00 Y.10 Y.01 Y.11
#>  0.5  0.5  0.5  0.5  0.5  0.5
```

You can also add custom priors. Custom priors are most simply specified by being added as a vector of numbers using **set_priors**. For instance:

```
R> make_model("X -> Y") |>
+ set_priors(1:6) |>
+ get_priors()

#>  X.0  X.1 Y.00 Y.10 Y.01 Y.11
#>   1    2    3    4    5    6
```

The priors here should be interpreted as indicating:

- $\alpha_X = (1, 2)$, which implies a distribution over $(\lambda_0^X, \lambda_1^X)$ centered on $(1/3, 2/3)$.
- $\alpha_Y = (3, 4, 5, 6)$, which implies a distribution over $(\lambda_{00}^Y, \lambda_{10}^Y, \lambda_{01}^Y, \lambda_{11}^Y)$ centered on $(3/18, 4/18, 5/18, 6/18)$.

For larger models it can be hard to provide priors as a vector of numbers. For that reason `set_priors` allows for more targeted modifications of the parameter vector. For instance:

```
R> make_model("X -> Y") |>
+ set_priors(statement = "Y[X=1] > Y[X=0]", alphas = 3) |>
+ get_priors()

#>  X.0  X.1 Y.00 Y.10 Y.01 Y.11
#>   1    1    1    1    3    1
```

As setting priors simply requires mapping alpha values to parameters, the process boils down to selecting rows of the `parameters_df` data frame, at which to alter values. When specifying a causal statement as above, *CausalQueries* internally identifies nodal types that are consistent with the statement, which in turn identify parameters to alter priors for.

Note that there is nothing particularly unique about setting priors via causal syntax, as it is simply an abstract way of specifying how to subset the `parameters_df` data frame. We can achieve the same result as above by specifying nodal types for which we would like to adjust the priors:

```
R> make_model("X -> Y") |>
+ set_priors(nodal_type = "01", alphas = 3) |>
+ get_priors()

#>  X.0  X.1 Y.00 Y.10 Y.01 Y.11
#>   1    1    1    1    3    1
```

or even parameter names

```
R> make_model("X -> Y") |>
+ set_priors(param_names = "Y.01", alphas = 3) |>
+ get_priors()

#>  X.0  X.1 Y.00 Y.10 Y.01 Y.11
#>   1    1    1    1    3    1
```

As such `set_priors` allows for the specification of any non-redundant combination of arguments on the `param_names`, `node`, `nodal_type`, `param_set` and `given` columns of `parameters_df` to uniquely identify parameters to set priors for. Alternatively a fully formed

subsetting statement may be supplied to `alter_at`. Since all these arguments get mapped to the parameters they identify internally they may be used interchangeably.⁴

```
R> make_model("X -> M -> Y; X <-> Y") |>
+ set_priors(node = "Y",
+           nodal_type = c("01","11"),
+           given = "X.1",
+           alphas = c(3,2)) |>
+ get_priors()

#>      X.0      X.1      M.00      M.10      M.01      M.11 Y.00_X.0 Y.10_X.0
#>      1      1      1      1      1      1      1      1
#> Y.01_X.0 Y.11_X.0 Y.00_X.1 Y.10_X.1 Y.01_X.1 Y.11_X.1
#>      1      1      1      1      3      2

R> make_model("X -> M -> Y; X <-> Y") |>
+ set_priors(
+   alter_at =
+     "node == 'Y' & nodal_type %in% c('01','11') & given == 'X.1'",
+   alphas = c(3,2)) |>
+ get_priors()

#>      X.0      X.1      M.00      M.10      M.01      M.11 Y.00_X.0 Y.10_X.0
#>      1      1      1      1      1      1      1      1
#> Y.01_X.0 Y.11_X.0 Y.00_X.1 Y.10_X.1 Y.01_X.1 Y.11_X.1
#>      1      1      1      1      3      2
```

It should be noted that while highly targeted prior setting is convenient and flexible; it should be done with caution. Setting priors on specific parameters in complex models; especially models involving confounding, may strongly affect inferences in intractable ways.

We additionally note that flat priors over nodal types do not necessarily translate into flat priors over queries. “Flat” priors over parameters in a parameter family put equal weight on each nodal type, but this in turn can translate into strong assumptions on causal quantities of interest.

For instance in an $X \rightarrow Y$ model in which negative effects are ruled out, the average causal effect implied by “flat” priors is $1/3$. This can be seen by querying the model as follows:

```
R> make_model("X -> Y") |>
+ set_restrictions(decreasing("X", "Y")) |>
+ query_model("Y[X=1] - Y[X=0]", n_draws = 10000)
```

More subtly the *structure* of a model, coupled with flat priors, has substantive importance for priors on causal quantities. For instance with flat priors, priors on the probability that X has a positive effect on Y in the model $X \rightarrow Y$ is centered on $1/4$. But priors on the probability that X has a positive effect on Y in the model $X \rightarrow M \rightarrow Y$ is centered on $1/8$.

Again, you can use `query_model` to figure out what flat (or other) priors over parameters imply for priors over causal quantities:

⁴See `?set_priors` and `?make_priors` for many more examples.

```
R> make_model("X -> Y") |>
+ query_model("Y[X=1] > Y[X=0]", n_draws = 10000)
R>
R> make_model("X -> M -> Y") |>
+ query_model("Y[X=1] > Y[X=0]", n_draws = 10000)
```

Caution regarding priors is particularly important when models are not identified, as is the case for many of the models considered here. In such cases, for some quantities, the marginal posterior distribution can be the same as the marginal prior distribution (Poirier 1998).

The key point here is to make sure you do not fall into a trap of thinking that “uninformative” priors make no commitments regarding the values of causal quantities of interest. They do, and the implications of flat priors for causal quantities can depend on the structure of the model. Moreover for some inferences from causal models the priors can matter a lot even if you have a lot of data. In such cases it can be helpful to know what priors on parameters imply for priors on causal quantities of interest (by using `query_model`) and to assess how much conclusions depend on priors (by comparing results across models that vary in their priors).

Setting Parameters

By default, models have a vector of parameter values included in the `parameters_df` dataframe. These are useful for generating data, or for situations, such as process tracing, when one wants to make inferences about causal types (θ), given case level data, under the assumption that the model is known.

The logic for setting parameters is similar to that for setting priors: effectively we need to place values on the probability of nodal types. The key difference is that whereas the *alpha* value placed on a nodal types can be any positive number—capturing our certainty over the parameter value—the parameter values must lie in the unit interval, $[0, 1]$. In general if parameter values are passed that do not lie in the unit interval, these are normalized so that they do.

Consider the causal model below. It has two parameter sets, one for X and one for Y , with six nodal types, two corresponding to X and four corresponding to Y . The key feature of the parameters is that they must sum to 1 within each parameter set.

```
R> make_model("X -> Y") |>
+ get_parameters()

#>   X.0  X.1 Y.00 Y.10 Y.01 Y.11
#> 0.50 0.50 0.25 0.25 0.25 0.25
```

The example below illustrates a change in the value of the parameter Y in the case it is increasing in X . Here nodal type $Y.Y01$ is set to be 0.5, while the other nodal types of this parameter set were renormalized so that the parameters in the set still sum to one.

```
R> make_model("X -> Y") |>
+ set_parameters(statement = "Y[X=1] > Y[X=0]", parameters = .5) |>
+ get_parameters()
```

```
#>      X.0      X.1      Y.00      Y.10      Y.01      Y.11
#> 0.5000000 0.5000000 0.1666667 0.1666667 0.5000000 0.1666667
```

5.4. Drawing and manipulating data

Once a model has been defined it is possible to simulate data from the model using the `make_data` function. This can be useful for instance for assessing the expected performance of a model given data drawn from some speculated set of parameter values.

```
R> model <- make_model("X -> M -> Y")
```

Drawing data basics

By default, the parameters used are taken from `model$parameters_df`.

```
R> sample_data_1 <-
+   model |>
+   make_data(n = 4)
```

However you can also specify parameters directly or use parameter draws from a prior or posterior distribution. For instance:

```
R> make_data(model, n = 3, param_type = "prior_draw")

#>   X M Y
#> 1 1 0 0
#> 2 1 1 0
#> 3 1 1 0
```

Note that the data is returned ordered by data type as in the example above.

Drawing incomplete data

`CausalQueries` can be used in settings in which researchers have gathered different amounts of data for different nodes. For instance gathering X and Y data for all units but M data only for some.

The function `make_data` allows you to draw data like this if you specify a data strategy indicating the probabilities of observing data on different nodes, possibly as a function of prior nodes observed.

```
R> set.seed(1)
R>
R> sample_data_2 <-
+   make_data(model,
+             n = 8,
+             nodes = list(c("X", "Y"), "M"),
+             probs = list(1, .5),
+             subsets = list(TRUE, "X==1 & Y==0"))
```

```
#> # A tibble: 2 x 5
#>   node_names nodes      n_steps probs subsets
#>   <chr>      <list>    <lgl>   <dbl> <chr>
#> 1 X, Y      <chr [2]> NA        1    TRUE
#> 2 M         <chr [1]> NA        0.5 X==1 & Y==0

R> sample_data_2

#>   X M Y
#> 1 0 NA 0
#> 2 0 NA 1
#> 3 1 NA 0
#> 4 1 NA 1
#> 5 1 NA 1
#> 6 1  1 0
#> 7 1 NA 0
#> 8 1  1 0
```

Reshaping data

Whereas data naturally comes in long form, with a row per observation, as in the examples above, the data passed to `stan` is in a compact form, which records only the number of units of each data type. `CausalQueries` includes functions that lets you move between these two forms in case of need.

```
R> sample_data_1 |> collapse_data(model)

#>   event strategy count
#> 1 XOM0Y0      XMY      0
#> 2 X1M0Y0      XMY      0
#> 3 XOM1Y0      XMY      1
#> 4 X1M1Y0      XMY      0
#> 5 XOM0Y1      XMY      0
#> 6 X1M0Y1      XMY      0
#> 7 XOM1Y1      XMY      3
#> 8 X1M1Y1      XMY      0

R> sample_data_2 |> collapse_data(model)

#>   event strategy count
#> 1 XOM0Y0      XMY      0
#> 2 X1M0Y0      XMY      0
#> 3 XOM1Y0      XMY      0
#> 4 X1M1Y0      XMY      2
#> 5 XOM0Y1      XMY      0
#> 6 X1M0Y1      XMY      0
#> 7 XOM1Y1      XMY      0
#> 8 X1M1Y1      XMY      0
#> 9 XOY0        XY       1
```



```
#> 10  X1Y0      XY      2
#> 11  X0Y1      XY      1
#> 12  X1Y1      XY      2
```

In the same way it is possible to move from “compact data” to “long data” using `expand_data()`. Note that NA’s are interpreted as data not having been sought. So in the case of `sample_data_2` the interpretation is that there are two data strategies: data on Y , M and X was sought in two cases only; data on Y and X only was sought in six cases.

6. Updating models

The approach used by the `CausalQueries` package to updating parameter values given observed data uses Stan ([Carpenter, Gelman, Hoffman, Lee, Goodrich, Betancourt, Brubaker, Guo, Li, and Riddell 2017](#)).

Below we explain the data required by the generic Stan program implemented in the package, the structure of that program, and then show how to use the package to produce posterior draws of parameters.

6.1. Data for stan

We use a generic Stan program that works for all binary causal models. The main advantage of the generic program we implement is that it allows us to pass the details of causal model as data inputs to Stan instead of generating individual Stan program for each causal model.

The data required by the Stan program includes vectors of observed data (Y) and priors on parameters (`lambdas_prior`) as well as a set of matrices required for the mapping between events, data types, causal types and parameters. The latter includes:

- A parameter matrix (P) that tells Stan how many parameters there are, and how they map into causal types,
- A matrix that maps parameters to data types (`parmap`), and
- An event matrix (E) that relates data types into patterns of observed data (events) in cases where there are incomplete observations.

In addition data includes counts of all relevant quantities as well as start and end positions of parameters pertaining to specific nodes and of distinct data strategies.

The internal function `prep_stan_data()` takes model and data as arguments and produces a list with all objects described above that are required by the generic Stan program. Generally, package users do not need to call the `prep_stan_data()` function directly to update the model. If further inspection of the data required by the Stan program is required, you can do so using the code below

```
R> sample_data_2 |>
+ collapse_data(model = model) |>
+ CausalQueries:::prep_stan_data(model = model)
```

6.2. How the Stan program works

The **stan** model involves the following elements:

- A specification of priors over sets of parameters
- A mapping from parameters to event probabilities, w
- A likelihood function.

Probability distributions over parameter sets

We are interested in “sets” of parameters. In the case without confounding these sets correspond to the nodal types for each node: we have a probability distribution over the set of nodal types. In cases with confounding these are sets of nodal types for a give node *given* values of other nodes: we have characterize the probability of each nodal type in a set given the values of nodal types for other nodes.

We express priors over these parameter sets using multiple Dirichlet distributions. In practice because we are dealing with multiple simplices of varying length, it is easier to express priors over gamma distributions with a unit scale parameter and shape parameter corresponding to the Dirichlet priors, α . We make use of the fact that $\lambda_0^X \sim \text{Gamma}(\alpha_0^X, 1)$ and $\lambda_1^X \sim \text{Gamma}(\alpha_1^X, 1)$ then $\frac{1}{\lambda_0^X + \lambda_1^X}(\lambda_0^X, \lambda_1^X) \sim \text{Dirichlet}(\alpha_0^X, \alpha_1^X)$ ⁵

TO illustrate, in the $X \rightarrow Y$ model we have two parameter sets (**param_sets**). The first is $\lambda^X \in \{\lambda_0^X, \lambda_1^X\}$ whose elements give the probability that X is 0 or 1. These two probabilities sum to one. The second parameter set is $\lambda^Y \in \{\lambda_{00}^Y, \lambda_{10}^Y, \lambda_{01}^Y, \lambda_{11}^Y\}$. These are also probabilities and their values sum to one. Note in all that we have 6 parameters but just $1 + 3 = 4$ degrees of freedom. We have a 2 dimensional Dirichlet distribution over the X nodal types (equivalently, a Beta distribution) and a 4 dimensional Dirichlet over the Y nodal types.

Event probabilities

For any candidate parameter vector λ we calculate the probability of *data* types. This is done using a matrix that maps from parameters into data types, **parmap**. In cases without confounding there is a column for each data type; the matrix indicates which nodes in each set “contribute” to the data type, and the probability of the data type is found by summing within sets and taking the product over sets.

```
R> make_model("X -> Y") |> get_parmap()
```

For example in table Figure 3 the probability of data type **X0Y0**, w_{00} is $\lambda_0^X \times \lambda_{00}^Y + \lambda_0^X \times \lambda_{01}^Y$. This is found by combining a parameter vector with the first column of **parmap**, taking the product of the probability of **X.0** and the *sum* of the probabilities for **Y.00** and **Y.01**.

In cases with confounding the approach is similar except that the **parmap** matrix can contain multiple columns for each data type to capture non-independence between nodes.

In the case of incomplete data we first identify the set of “data strategies”, where a collection of a data strategy might be of the form “gather data on X and M , but not Y , for n_1 cases

⁵For a discussion of implementation of this approach in **stan** see discussion (here)[<https://discourse.mc-stan.org/t/ragged-array-of-simplexes/1382>].

	X0Y0	X1Y0	X0Y1	X1Y1
X.0	1	0	1	0
X.1	0	1	0	1
Y.00	1	1	0	0
Y.10	0	1	1	0
Y.01	1	0	0	1
Y.11	0	0	1	1

Figure 3: ?(caption)

and gather data on X and Y , but not M , for n_2 cases. The probability of an observed event, within a data strategy, is given by summing the probabilities of the types that could give rise to the incomplete data. For example X is observed, but Y is not, then the probability of $X = 0, Y = \text{NA}$ is $w_{00} + w_{01}$. The matrix E is passed to **stan** to figure out which event probabilities need to be combined for events with missing data.

Data probability

Once we have the event probabilities in hand for each data strategy we are ready to calculate the probability of the data. For a given data strategy this is given by a multinomial distribution with these event probabilities. When there is incomplete data, and so multiple data strategies, this is given by the the product of the multinomial probabilities for each strategy.

6.3. Implementation

To update a **CausalQueries** model with data use:

```
R> update_model(model, data)
```

where the data argument is a dataset containing some or all of the nodes in the model. `update_model()` relies on `rstan::sampling()` to draw from posterior distribution and one can pass any additional arguments accepted by `rstan::sampling()` in

6.4. Parallelization

If you have multiple cores you can do parallel processing by including this line before running **CausalQueries**:

```
R> library(parallel)
R>
R> options(mc.cores = parallel::detectCores())
```

Additionally parallelizing across models or data while running MCMC chains in parallel can be achieved by setting up a nested parallel process. With 8 cores one can run 2 updating processes with 3 parallel chains each simultaneously. More generally the number of parallel processes at the upper level of the nested parallel structure are given by $\left\lfloor \frac{\text{cores}}{\text{chains} + 1} \right\rfloor$.

```

R> library(future)
R> library(future.apply)
R>
R> chains <- 3
R> cores <- 8
R>
R> future::plan(list(
+   future::tweak(future::multisession,
+                 workers = floor(cores/(chains + 1))),
+   future::tweak(future::multisession,
+                 workers = chains)
+ ))
R>
R> model <- make_model("X -> Y")
R> data <- list(data_1, data_2)
R>
R> future.apply::future_lapply(data, function(d) {
+   update_model(
+     model = model,
+     data = d,
+     chains = chains,
+     refresh = 0
+   )
+ })

```

6.5. Incomplete and censored data

CausalQueries assumes that missing data is missing at random, conditional on observed data. Thus for instance in a $X \rightarrow M \rightarrow Y$ model one might choose to observe M in a random set of cases in which $X = 1$ and $Y = 1$. In that case if there are positive relations at each stage you may be more likely to observe M in cases in which $M = 1$. However observation of M is still random *conditional* on the observed X and Y data. The **stan** model in *CausalQueries* takes account of this kind of sampling naturally by assessing the probability of observing a particular pattern of data within each data strategy. For a discussion see section 9.2.3.2 of [Humphreys and Jacobs \(2023\)](#).

In addition it is possible to indicate when data has been censored and for the **stan** model to take this into account also. Say for instance that we only get to observe X in cases where $X = 1$ and not when $X = 0$. This kind of sampling is non random conditional on observables. It is taken account however by indicating to **stan** that the probability of observing a particular data type is 0, regardless of parameter values. This is done using the `censored_types` argument in `update_model()`.

To illustrate, in the example below we observe perfectly correlated data for X and Y . If we are aware that data in which $X \neq Y$ has been censored then when we update we do not move towards a belief that X causes Y .

```

R> data <- data.frame(X = rep(0:1, 5), Y = rep(0:1, 5))

```

Table 10: ?(caption)

Table 11: Posterior inferences taking account of censoring and not.

model	query	mean	sd
uncensored	$(Y[X=1] - Y[X=0])$	0.59	0.20
censored	$(Y[X=1] - Y[X=0])$	0.01	0.31

```

R>
R> list(
+   uncensored = make_model("X -> Y") |>
+     update_model(data),
+
+   censored = make_model("X -> Y") |>
+     update_model(data, censored_types = c("X1Y0", "X0Y1"))
+) %>%
+   query_model(te("X", "Y"), using = "posteriors")

```

6.6. Output

The primary output from `update_model()` is a model with an attached posterior distribution over model parameters, stored as a dataframe in `model$posterior_distribution`. In addition, a distribution of causal types is stored by default and the `stanfit` object and a distribution over event probabilities are optionally saved.

7. Queries

7.1. Calculating factual and counterfactual quantities

We can map from causal types to data types by propagating realized values on nodes forward in the DAG, moving from exogenous or intervened upon nodes to their descendants in generational order. The `realise_outcomes` function achieves this by traversing the DAG, while recording for each node's nodal types, the values implied by realizations on the node's parents. By way of example, consider the first causal type of a $X \rightarrow Y$ model:

1. X is exogenous and has a realized value of 0
2. We substitute for Y the value implied by the 00 nodal type given a 0 value on X , which in turn is 0 (see nodal types)

Recovering implied values on complex nodal types efficiently at scale, exploits the fact that nodal types are the Cartesian Product of possible parent realizations. Finding the index of a Node's realized value in a nodal type given parent realizations is equivalent to finding the row index in the Cartesian Product matrix corresponding to those parent realizations. By definition of the Cartesian product, the number of consecutive 0 or 1 elements in a given

column is $2^{\text{columnindex}}$, when indexing columns from 0. Given a set of parent realizations R indexed from 0, the corresponding row index in the Cartesian Product Matrix indexed from 0 can thus be computed via: $\text{row} = (\sum_{i=0}^{|R|-1} (2^i \times R_i))$.

Calling `realise_outcomes` on the above model thus yields:

```
R> make_model("X -> Y") |> realise_outcomes()

#>      X Y
#> 0.00 0 0
#> 1.00 1 0
#> 0.10 0 1
#> 1.10 1 0
#> 0.01 0 0
#> 1.01 1 1
#> 0.11 0 1
#> 1.11 1 1
```

with row names indicating nodal types and columns realized values. Intervening on X with $\text{do}(X = 1)$ yields:

```
R> make_model("X -> Y") |> realise_outcomes(dos = list(X = 1))

#>      X Y
#> 0.00 1 0
#> 1.00 1 0
#> 0.10 1 0
#> 1.10 1 0
#> 0.01 1 1
#> 1.01 1 1
#> 0.11 1 1
#> 1.11 1 1
```

7.2. Causal Syntax

CausalQueries provides syntax for the formulation of various causal queries including queries on all rungs of the “causal ladder” (Pearl and Mackenzie 2018): prediction, such as the proportion of units where Y equals 1; intervention, such as the probability that $Y = 1$ when X is *set* to 1; counterfactuals, such as the probability that Y would be 1 were $X = 1$ given we know Y is 0 when X was observed to be 0. Queries can be posed at the population level or case level and can be unconditional (e.g. what is the effect of X on Y for all units) or conditional (for example, the effect of X on Y for units for whom Z affects X).

This syntax enables users to write arbitrary causal queries to interrogate their models. The heart of querying is figuring out which causal types correspond to particular queries.

For factual queries, users may employ logical statements to ask questions about observed conditions, without any intervention. Take, for example, the query mentioned above about the proportion of units where Y equals 1, expressed as `"Y == 1"`. In this case the logical operator `==` indicates that *CausalQueries* should consider units that fulfill the condition of

strict equality where Y equals 1.⁶ When this query is posed, the `get_query_types` function identifies all types that give rise to $Y = 1$, absent any interventions.

```
R> make_model("X -> Y") |>
+   get_query_types("Y==1")

#>
#> Causal types satisfying query's condition(s)
#>
#>   query =   Y==1
#>
#> X0.Y10  X1.Y01
#> X0.Y11  X1.Y11
#>
#>
#>   Number of causal types that meet condition(s) =   4
#>   Total number of causal types in model =   8
```

The key to posing causal queries is being able to ask about values of variables given that the values of some other variables are “controlled”. This corresponds to application of the `do` operator in Pearl (2009). In `CausalQueries` this is done by using square brackets `[]`. Inside the brackets, variables that are to be intervened upon are specified.

For instance, consider the query `Y[X=0]==1`. This query asks about the types for which Y equals 1 when X is *set* to 0. In this case, since X is being intervened to be zero, X is placed inside the brackets. Given that Y equaling 1 is a condition about potentially observed values, it is expressed as using the logical operator `==`.

The set of causal types that meets this query is quite different:

```
R> make_model("X -> Y") |>
+   get_query_types("Y[X=1]==1")

#>
#> Causal types satisfying query's condition(s)
#>
#>   query =   Y[X=1]==1
#>
#> X0.Y01  X1.Y01
#> X0.Y11  X1.Y11
#>
#>
#>   Number of causal types that meet condition(s) =   4
#>   Total number of causal types in model =   8
```

When a node has multiple parents it is possible to set the value of non of the parents, of some of the parents, or of all of the parents. For instance if X_1 and X_2 are parents of Y then `Y==1`, `Y[X1 = 1]==1`, and `Y[X1 = 1, X2 = 1]==1` queries cases for which $Y = 1$ when, respectively,

⁶`CausalQueries` also accepts `=` as a shorthand for `==`. However, `==` is preferred as it is the conventional logical operator to express a condition of strict equality.

neither parents values are controlled, when X_1 is set to 1 but X_2 is not controlled, and when both X_1 and X_2 are set to 1.

As an example we have:

```
R> make_model("X1 -> Y <- X2") |>
+ get_query_types("X1 ==1 & X2 == 1 & (Y[X1=1, X2=1] > Y[X1=0, X2=0])")

#>
#> Causal types satisfying query's condition(s)
#>
#> query = X1==1&X2==1&(Y[X1=1,X2=1]>Y[X1=0,X2=0])
#>
#> X11.X21.Y0001 X11.X21.Y0101
#> X11.X21.Y0011 X11.X21.Y0111
#>
#>
#> Number of causal types that meet condition(s) = 4
#> Total number of causal types in model = 64
```

In this case, the aim is to identify the types for which $X_1 = X_2 = 1$ and Y equals zero when $X_1 = 0$ and $X_2 = 0$, and Y equals 1 when $X_1 = 1$ and $X_2 = 1$.

```
R> make_model('X -> M -> Y <- X') |>
+ get_query_types(query = "Y[X=1, M=1] > Y[X=0, M=0]",
+ map = "nodal_type")

#>
#> Nodal types adding weight to query
#>
#> query : Y[X=1,M=1]>Y[X=0,M=0]
#>
#> 0001 0101
#> 0011 0111
#>
#>
#> Number of nodal types that add weight to query = 4
#> Total number of nodal types related to Y = 16
```

In summary, the variables to be intervened, as if conducting an experiment, are placed inside square brackets, followed by an equal sign and the value to which we want to set them, either 1 or 0, as X in the case of " $Y[X=1]$ ". The variable whose value is to be observed, as Y in the same example, should be placed before the square brackets. Finally, conditions related to observed or potentially observed values, in the context of an intervention, are expressed outside the brackets, along with the logical condition that defines the observed values, as in " $Y == 1$ " or " $Y[X=1] > Y[X=0]$ ".

Complex expressions

Many queries involve complex statements over multiple sets of types. These can be formed

with the aid of logical operators such as `==`, `!=`, `>`, `>=`, `<`, `<=`. For example, they can make queries about cases where X has a positive effect on Y , i.e., whether Y is greater when X is set to 1 compared to when X is set to 0, expressed as `"Y[X=1] > Y[X=0]"`. The query “ X has some effect on Y ” is given by `"Y[X=1] != Y[X=0]"`.

Linear operators can also be used over set of simple statements. Thus `"Y[X=1] - Y[X=0]"` returns the share of units the average treatment effect. In essence rather than returning a `TRUE` or `FALSE` for the two parts of the query, the case memberships are forced to numeric values (1 or 0) and the differences taken, which can be a 1, 0 or -1 depending on the causal type.

```
R> make_model("X -> Y") |>
+   get_query_types("Y[X=1]- Y[X=0]")

#> X0.Y00 X1.Y00 X0.Y10 X1.Y10 X0.Y01 X1.Y01 X0.Y11 X1.Y11
#>      0      0      -1      -1      1      1      0      0
```

Nested queries

`CausalQueries` lets users pose nested “complex counterfactual” queries. For instance `"Y[M=M[X=0], X=1]==1"` queries the types for which Y equals 1 when X is set to 1, while keeping M constant at the value it would take if X were 0.

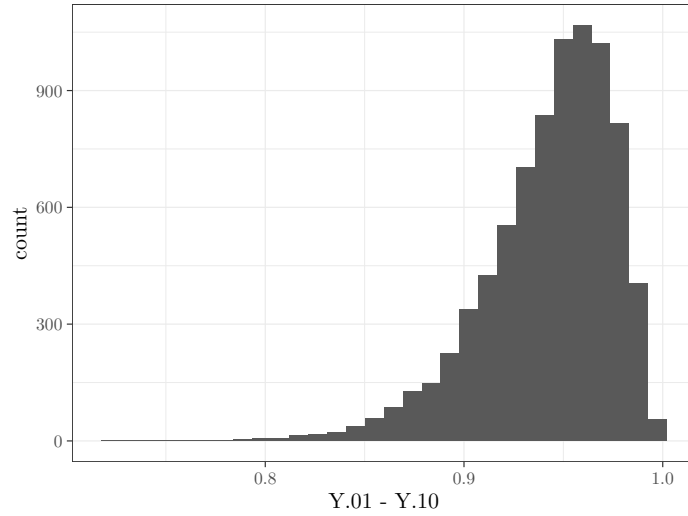
7.3. Quantifying queries

Giving a *quantitative* answer to a query requires placing probabilities over the causal types that correspond to a query.

Queries by hand

Queries can be calculated directly from the prior distribution of the posterior distribution provided by `stan`.

```
R> data <- data.frame(X = rep(0:1, 50), Y = rep(0:1, 50))
R>
R> model <-
+   make_model("X -> Y") |>
+   update_model(data, iter = 4000)
R>
R> model$posterior_distribution |>
+   ggplot(aes(Y.01 - Y.10)) + geom_histogram()
```

Figure 4: Posterior on 'Probability Y is increasing in X '*Query distribution*

It is generally useful to use causal syntax to define the query and calculate the query with respect to the prior or posterior probability distributions.

This can be done for a list of queries using `query_distribution`:

```
R> make_model("X -> Y") |>
+ query_distribution(
+   query = list(increasing = "(Y[X=1] > Y[X=0])"),
+   using = "priors") |>
+ ggplot(aes(increasing)) + geom_histogram() + theme_bw()
```

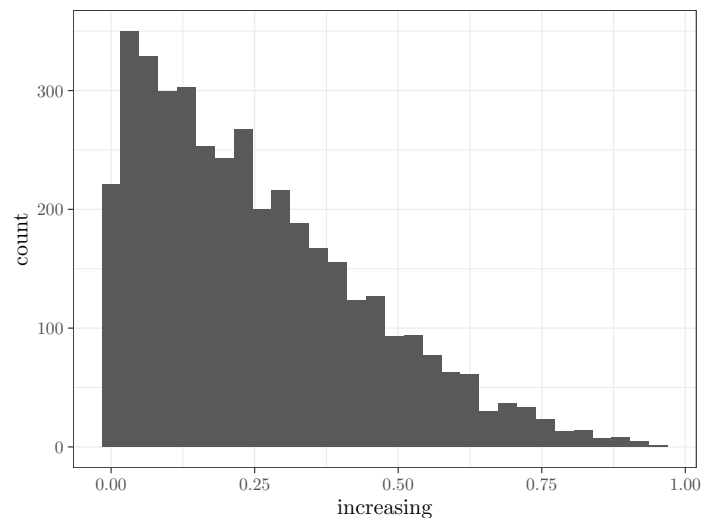
`query_distribution` can also be used when one is interested in assessing the value of a query for a *particular case*. In a sense this is equivalent to posing a conditional query, querying conditional on values in a case. For instance we might consult our posterior lipids model and ask about the effect of X on Y for a case in which $Z = 1$, $X = 1$ and $Y = 1$.

```
R> lipids_model |>
+ query_model(query = "Y[X=1] - Y[X=0]",
+   given = c("X==1 & Y==1 & Z==1"),
+   using = "posteriors")
```

Table 12: Case Level Query Example.

query	given	mean	sd	cred.low.2.5%	cred.high.97.5%
$Y[X=1] - Y[X=0]$	$X==1 \ \& \ Y==1 \ \& \ Z==1$	0.95	0.04	0.87	1

The answer we get in Table 12 is what we now believe for all cases in which $Z = 1$, $X = 1$ and $Y = 1$. It is in fact the expected average effect among cases with this data type and so this expectation has an uncertainty attached to it.

Figure 5: Prior on 'Probability Y is increasing in X '

Subtly though this is, in principle, different to what we would infer for a “new case” that we wonder about. When inquiring about a new case, the case level query *updates* on the given information observed in the new case. The resulting inference is different to the inference that would be made from the posterior *given* the features of the case.

To illustrate, consider a model for which we are quite sure that X causes Y but we do not know whether it works through two positive effects or two negative effects.

Thus we do not know if $M=0$ would suggest an effect or no effect. If asked what we would infer for a case that had $M = 0$ we would not know whether $M = 0$ information is consistent with a positive effect or not. However if provided with a randomly case and learn that it has $M = 0$, then we update about the causal model and infer that there is an effect in this case (but that there wouldn't be were $M = 1$).

```
R> # DO ALL THIS IN ONE TABLE
R>
R> set.seed(1)
R> model <-
+   make_model("X -> M -> Y") |>
+   update_model(data.frame(X = rep(0:1, 8), Y = rep(0:1, 8)), iter = 10000)
R>
R> Q <- "Y[X=1] > Y[X=0]"
R> G <- "X==1 & Y==1 & M==1"
R> QG <- "(Y[X=1] > Y[X=0]) & (X==1 & Y==1 & M==1)"
R>
R> # In this case these are very different:
R> query_distribution(model, Q, given = G, using = "posteriors")[[1]] |> mean()
R> query_distribution(model, Q, given = G, using = "posteriors",
+   case_level = TRUE)
R>
```

```

R> # These are equivalent:
R> # 1. Case level query via function
R> query_distribution(model, Q, given = G,
+   using = "posteriors", case_level = TRUE)
R>
R> # 2. Case level query by hand using Bayes
R> distribution <- query_distribution(
+   model, list(QG = QG, G = G), using = "posteriors")
R>
R> mean(distribution$QG)/mean(distribution$G)

```

Batch queries

The `query_model()` function takes list of models, causal queries, and conditions (**given**), calculates estimands given prior or posterior distributions and reports summaries of these distributions. The result is a data frame which can be displayed as a table.

```

R> models <- list(
+ `1` = make_model("X -> Y"),
+ `2` = make_model("X -> Y") |> set_restrictions("Y[X=1] < Y[X=0]")
+ )
R>
R> query_model(
+ models,
+ query = list(ATE = "Y[X=1] - Y[X=0]",
+               POS = "Y[X=1] > Y[X=0]"),
+ given = c(TRUE, "Y==1 & X==1"),
+ case_level = c(FALSE, TRUE),
+ using = c("parameters", "priors"),
+ expand_grid = TRUE)

```

Table 13: Results for Two Queries on Two Models.

model	query	given	using	case_level	mean	sd
1	ATE	-	parameters	FALSE	0.00	NA
2	ATE	-	parameters	FALSE	0.33	NA
1	ATE	-	priors	FALSE	-0.01	0.32
2	ATE	-	priors	FALSE	0.33	0.24
1	ATE	Y==1 & X==1	parameters	FALSE	0.50	NA
2	ATE	Y==1 & X==1	parameters	FALSE	0.50	NA
1	ATE	Y==1 & X==1	priors	FALSE	0.50	0.29
2	ATE	Y==1 & X==1	priors	FALSE	0.50	0.29
1	POS	-	parameters	FALSE	0.25	NA
2	POS	-	parameters	FALSE	0.33	NA
1	POS	-	priors	FALSE	0.25	0.20
2	POS	-	priors	FALSE	0.33	0.24
1	POS	Y==1 & X==1	parameters	FALSE	0.50	NA

2	POS	Y==1 & X==1	parameters	FALSE	0.50	NA
1	POS	Y==1 & X==1	priors	FALSE	0.50	0.29
2	POS	Y==1 & X==1	priors	FALSE	0.50	0.29
1	ATE	-	parameters	TRUE	0.00	NA
2	ATE	-	parameters	TRUE	0.33	NA
1	ATE	-	priors	TRUE	-0.01	NA
2	ATE	-	priors	TRUE	0.33	NA
1	ATE	Y==1 & X==1	parameters	TRUE	0.50	NA
2	ATE	Y==1 & X==1	parameters	TRUE	0.50	NA
1	ATE	Y==1 & X==1	priors	TRUE	0.51	NA
2	ATE	Y==1 & X==1	priors	TRUE	0.50	NA
1	POS	-	parameters	TRUE	0.25	NA
2	POS	-	parameters	TRUE	0.33	NA
1	POS	-	priors	TRUE	0.25	NA
2	POS	-	priors	TRUE	0.33	NA
1	POS	Y==1 & X==1	parameters	TRUE	0.50	NA
2	POS	Y==1 & X==1	parameters	TRUE	0.50	NA
1	POS	Y==1 & X==1	priors	TRUE	0.51	NA
2	POS	Y==1 & X==1	priors	TRUE	0.50	NA

The `stats` argument let users specify alternative summary statistics and indicate. The `expand_grid` argument lets users specify whether all combinations of list elements should be generated.

Computational details and software requirements

Version	<ul style="list-style-type: none"> 1.0.1
Availability	<ul style="list-style-type: none"> Stable <p>Release:</p> <p>https://cran.rstudio.com/web/packages/CausalQueries/index.html</p> <ul style="list-style-type: none"> Development: <p>https://github.com/integrated-inferences/CausalQueries</p>
Issues	<ul style="list-style-type: none"> https://github.com/integrated-inferences/CausalQueries/issues
Operating Systems	<ul style="list-style-type: none"> Linux MacOS Windows
Testing Environments	<ul style="list-style-type: none"> Ubuntu 22.04.2 Debian 12.2
OS	<ul style="list-style-type: none"> MacOS Windows
Testing Environments	<ul style="list-style-type: none"> R 4.3.1 R 4.3.0
R	<ul style="list-style-type: none"> R 4.2.3 r-devel

R Version	• R(>= 3.4.0)
Compiler	• either of the below or similar:
	• g++
	• clang++
Stan	• inline
requirements	• Rcpp (>= 0.12.0)
	• RcppEigen (>= 0.3.3.3.0)
	• RcppArmadillo (>= 0.12.6.4.0)
	• RcppParallel (>= 5.1.4)
	• BH (>= 1.66.0)
	• StanHeaders (>= 2.26.0)
	• rstan (>= 2.26.0)
R-Packages	• dplyr
Depends	• methods
R-Packages	• dagitty (>= 0.3-1)
Imports	• dirmult (>= 0.1.3-4)
	• stats (>= 4.1.1)
	• rlang (>= 0.2.0)
	• rstan (>= 2.26.0)
	• rstantools (>= 2.0.0)
	• stringr (>= 1.4.0)
	• ggdag (>= 0.2.4)
	• latex2exp (>= 0.9.4)
	• ggplot2 (>= 3.3.5)
	• lifecycle (>= 1.0.1)

The results in this paper were obtained using R~3.4.1 with the **MASS**~7.3.47 package. R itself and all packages used are available from the Comprehensive R Archive Network (CRAN) at <https://CRAN.R-project.org/>.

Acknowledgments

The approach to generating a generic `stan` function that can take data from arbitrary models was developed in key contributions by [Jasper Cooper](#) and [Georgiy Syunyaev](#). [Lily Medina](#) did magical work pulling it all together and developing approaches to characterizing confounding and defining estimands. Clara Bicalho helped figure out the syntax for causal statements. Julio Solis made many key contributions figuring out how to simplify the specification of priors. Merlin Heidemanns figured out the `rstantools` integration and made myriad code improvements. Till Tietz revamped the entire package and improved every part of it.

References

- Angrist JD, Imbens GW, Rubin DB (1996). “Identification of causal effects using instrumental variables.” *Journal of the American statistical Association*, **91**(434), 444–455.
- Balke A, Pearl J (1997). “Bounds on treatment effects from studies with imperfect compliance.” *Journal of the American Statistical Association*, **92**(439), 1171–1176.
- Beaumont P, Horsburgh B, Pilgerstorfer P, Droth A, Oentaryo R, Ler S, Nguyen H, Ferreira GA, Patel Z, Leong W (2021). “CausalNex.” URL <https://github.com/quantumblacklabs/causalnex>.
- Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M, Brubaker M, Guo J, Li P, Riddell A (2017). “Stan: A probabilistic programming language.” *Journal of statistical software*, **76**(1).
- Chickering DM, Pearl J (1996). “A clinician’s tool for analyzing non-compliance.” In *Proceedings of the National Conference on Artificial Intelligence*, pp. 1269–1276.
- Duarte G, Finkelstein N, Knox D, Mummolo J, Shpitser I (2023). “An Automated Approach to Causal Inference in Discrete Settings.” *Journal of the American Statistical Association*, **0**(0), 1–16. doi:10.1080/01621459.2023.2216909. <https://doi.org/10.1080/01621459.2023.2216909>, URL <https://doi.org/10.1080/01621459.2023.2216909>.
- Frangakis CE, Rubin DB (2002). “Principal stratification in causal inference.” *Biometrics*, **58**(1), 21–29.
- Humphreys M, Jacobs AM (2023). *Integrated Inferences*. Cambridge University Press.
- Kalisch M, Mächler M, Colombo D, Maathuis MH, Bühlmann P (2012). “Causal inference using graphical models with the R package pcalg.” *Journal of statistical software*, **47**, 1–26.
- Michael C Sachs, Gustav Jonzon AS, Gabriel EE (2023). “A General Method for Deriving Tight Symbolic Bounds on Causal Effects.” *Journal of Computational and Graphical Statistics*, **32**(2), 567–576. doi:10.1080/10618600.2022.2071905. <https://doi.org/10.1080/10618600.2022.2071905>, URL <https://doi.org/10.1080/10618600.2022.2071905>.
- Pearl J (2009). *Causality*. Cambridge university press.
- Pearl J, Mackenzie D (2018). *The book of why: the new science of cause and effect*. Basic Books.
- Poirier DJ (1998). “Revising beliefs in nonidentified models.” *Econometric Theory*, **14**(4), 483–509.
- Zhang J, Tian J, Bareinboim E (2022). “Partial Counterfactual Identification from Observational and Experimental Data.” In K Chaudhuri, S Jegelka, L Song, C Szepesvari, G Niu, S Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 26548–26558. PMLR. URL <https://proceedings.mlr.press/v162/zhang22ab.html>.

More technical details

Appendices can be included after the bibliography (with a page break). Each section within the appendix should have a proper section title (rather than just *Appendix*). For more technical style details, please check out JSS's style FAQ at [<https://www.jstatsoft.org/pages/view/style#frequently-asked-questions>] which includes the following topics:

- Title vs. sentence case.
- Graphics formatting.
- Naming conventions.
- Turning JSS manuscripts into R package vignettes.
- Trouble shooting.
- Many other potentially helpful details...

Using BibTeX

References need to be provided in a BibTeX file (`.bib`). All references should be made with `@cite` syntax. This commands yield different formats of author-year citations and allow to include additional details (e.g., pages, chapters, ...) in brackets. In case you are not familiar with these commands see the JSS style FAQ for details.

Cleaning up BibTeX files is a somewhat tedious task – especially when acquiring the entries automatically from mixed online sources. However, it is important that informations are complete and presented in a consistent style to avoid confusions. JSS requires the following format.

- item JSS-specific markup (`\proglang`, `\pkg`, `\code`) should be used in the references.
- item Titles should be in title case.
- item Journal titles should not be abbreviated and in title case.
- item DOIs should be included where available.
- item Software should be properly cited as well. For R packages `citation("pkgname")` typically provides a good starting point.

8. Appendix: stan code

`prep_stan_data` then returns a list of objects that `stan` expects to receive. These include indicators to figure out where a parameter set starts (`l_starts`, `l_ends`) and ends and where a data strategy starts and ends (`strategy_starts`, `strategy_ends`), as well as the matrices described above.

Below we show the `stan` code. This starts off with a block saying what input data is to be expected. Then there is a characterization of parameters and the transformed parameters. Then the likelihoods and priors are provided. `stan` takes it from there and generates a posterior distribution.


```

#> functions{
#>   row_vector col_sums(matrix X) {
#>     row_vector[cols(X)] s ;
#>     s = rep_row_vector(1, rows(X)) * X ;
#>     return s ;
#>   }
#> }
#> data {
#> int<lower=1> n_params;
#> int<lower=1> n_paths;
#> int<lower=1> n_types;
#> int<lower=1> n_param_sets;
#> int<lower=1> n_nodes;
#> array[n_param_sets] int<lower=1> n_param_each;
#> int<lower=1> n_data;
#> int<lower=1> n_events;
#> int<lower=1> n_strategies;
#> int<lower=0, upper=1> keep_transformed;
#> vector<lower=0>[n_params] lambdas_prior;
#> array[n_param_sets] int<lower=1> l_starts;
#> array[n_param_sets] int<lower=1> l_ends;
#> array[n_nodes] int<lower=1> node_starts;
#> array[n_nodes] int<lower=1> node_ends;
#> array[n_strategies] int<lower=1> strategy_starts;
#> array[n_strategies] int<lower=1> strategy_ends;
#> matrix[n_params, n_types] P;
#> matrix[n_params, n_paths] parmap;
#> matrix[n_paths, n_data] map;
#> matrix<lower=0, upper=1>[n_events, n_data] E;
#> array[n_events] int<lower=0> Y;
#> }
#> parameters {
#> vector<lower=0>[n_params - n_param_sets] gamma;
#> }
#> transformed parameters {
#> vector<lower=0, upper=1>[n_params] lambdas;
#> vector<lower=1>[n_param_sets] sum_gammas;
#> matrix[n_params, n_paths] parlam;
#> matrix[n_nodes, n_paths] parlam2;
#> vector<lower=0, upper=1>[n_paths] w_0;
#> vector<lower=0, upper=1>[n_data] w;
#> vector<lower=0, upper=1>[n_events] w_full;
#> // Cases in which a parameter set has only one value need special handling
#> // they have no gamma components and sum_gamma needs to be made manually
#> for (i in 1:n_param_sets) {
#>   if (l_starts[i] >= l_ends[i]) {
#>     sum_gammas[i] = 1;

```

```

#> // syntax here to return unity as a vector
#> lambdas[l_starts[i]] = lambdas_prior[1]/lambdas_prior[1];
#> }
#> else if (l_starts[i] < l_ends[i]) {
#>   sum_gammas[i] =
#>   1 + sum(gamma[(l_starts[i] - (i-1)):(l_ends[i] - i)]);
#>   lambdas[l_starts[i]:l_ends[i]] =
#>   append_row(1, gamma[(l_starts[i] - (i-1)):(l_ends[i] - i)]) /
#>   sum_gammas[i];
#> }
#> }
#> // Mapping from parameters to data types
#> // (usual case): [n_par * n_data] * [n_par * n_data]
#> parlam = rep_matrix(lambdas, n_paths) .* parmap;
#> // Sum probability over nodes on each path
#> for (i in 1:n_nodes) {
#>   parlam2[i,] = col_sums(parlam[(node_starts[i]):(node_ends[i]),]);
#> }
#> // then take product to get probability of data type on path
#> for (i in 1:n_paths) {
#>   w_0[i] = prod(parlam2[,i]);
#> }
#> // last (if confounding): map to n_data columns instead of n_paths
#> w = map'*w_0;
#> // Extend/reduce to cover all observed data types
#> w_full = E * w;
#> }
#> model {
#> // Dirichlet distributions (earlier versions used gamma)
#> for (i in 1:n_param_sets) {
#>   target += dirichlet_lpdf(lambdas[l_starts[i]:l_ends[i]] |
#>     lambdas_prior[l_starts[i] :l_ends[i]]);
#>   target += -n_param_each[i] * log(sum_gammas[i]);
#> }
#> // Multinomials
#> // Note with censoring event_probabilities might not sum to 1
#> for (i in 1:n_strategies) {
#>   target += multinomial_lpmf(
#>     Y[strategy_starts[i]:strategy_ends[i]] |
#>     w_full[strategy_starts[i]:strategy_ends[i]]/
#>     sum(w_full[strategy_starts[i]:strategy_ends[i]]));
#> }
#> }
#> // Option to export distribution of causal types
#> generated_quantities{
#> vector[n_types] prob_of_types;
#> if (keep_transformed == 1){

```

```
#> for (i in 1:n_types) {  
#>   prob_of_types[i] = prod(P[, i].*lambdas + 1 - P[,i]);  
#> }  
#> if (keep_transformed == 0){  
#>   prob_of_types = rep_vector(1, n_types);  
#> }  
#> }
```

Affiliation:

Till Tietz

IPI

Reichpietschufer 50

Berlin Germany

E-mail: ttietz2014@gmail.com

URL: <https://github.com/till-tietz>

Lily Medina

E-mail: lily.medina@berkeley.edu

URL: <https://lilymedina.github.io/>

Georgiy Syunyaev

E-mail: g.syunyaev@vanderbilt.edu

URL: <https://gsyunyaev.com/>

Macartan Humphreys

E-mail: macartan.humphreys@wzb.eu

URL: <https://macartan.github.io/>