

## Reproduction Information

### **Causal inference to scope environmental impact assessment of renewable energy projects and test competing mental models of decarbonization**

Amir M. Gazar<sup>1,2</sup>, Mark E. Borsuk<sup>3</sup>, and Ryan S.D. Calder<sup>1,2,3,4,5,\*</sup>

<sup>1</sup>Department of Population Health Sciences, Virginia Tech, Blacksburg, VA, 24061, USA

<sup>2</sup>Global Change Center, Virginia Tech, Blacksburg, VA, 24061, USA

<sup>3</sup>Department of Civil and Environmental Engineering, Duke University, Durham, NC, 27708, USA

<sup>4</sup>Faculty of Health Sciences, Virginia Tech, Roanoke, VA, 24016, USA

<sup>5</sup>Department of Civil and Environmental Engineering, Virginia Tech, Blacksburg, VA, 24061, USA

\* Corresponding author: [rsdc@vt.edu](mailto:rsdc@vt.edu)

# Contents

<b>1</b>	<b>Configuration and Setup</b>	<b>3</b>
1.1	Hardware and Operating System . . . . .	3
1.2	Software Versions and Dependencies . . . . .	3
1.3	Install Time . . . . .	3
1.4	Run Time . . . . .	3
<b>2</b>	<b>Full code execution instructions</b>	<b>3</b>
<b>3</b>	<b>Pseudocode Overview</b>	<b>4</b>
<b>4</b>	<b>Algorithm Overview</b>	<b>5</b>
<b>5</b>	<b>Preparing the R Environment</b>	<b>5</b>
5.1	Install and Load the Required Packages . . . . .	5
5.2	Set Working Directory . . . . .	6
<b>6</b>	<b>Required Functions</b>	<b>6</b>
6.1	DAG Visualizer . . . . .	6
6.2	Box-Cox Transformer . . . . .	6
6.3	Goodness of Fit Test . . . . .	7
6.4	D-Separation Test . . . . .	7
6.4.1	Function Usage . . . . .	8
6.4.2	Parameters . . . . .	8
6.4.3	Return Value . . . . .	8
6.4.4	Underpinning Functions . . . . .	9
<b>7</b>	<b>Data Preparation</b>	<b>10</b>
7.1	Generating delta and lags . . . . .	10
7.2	Creating averaged and lagged variables . . . . .	11
7.3	Creating a subset for averaged and lagged variables . . . . .	11
7.4	Ensuring data-set is numeric . . . . .	11
7.5	Creating a new data-frame . . . . .	11
7.6	Check for Gaussian distribution . . . . .	12
7.7	Transforming non-Gaussian variables . . . . .	12
7.8	Discretizing variables . . . . .	12
7.9	Creating the final data-frame . . . . .	12
<b>8</b>	<b>Creating a Blacklist</b>	<b>13</b>
<b>9</b>	<b>Visualizing the Hypothesized DAG</b>	<b>13</b>
<b>10</b>	<b>Constructing Score-Based DAGs</b>	<b>14</b>
<b>11</b>	<b>Plotting Conditional Dependency Results</b>	<b>15</b>
11.1	Running the Graph_Generator.R Code . . . . .	15
11.2	Snippet of Graph_Generator.R Code . . . . .	15
<b>12</b>	<b>Goodness of Fit for Each Node</b>	<b>16</b>
<b>13</b>	<b>D-Separation Analysis</b>	<b>16</b>

## 1 Configuration and Setup

This document presents a comprehensive guide outlining the sequential process required to execute and replicate the causal inference modeling methodology presented in our article. Additionally, we supplied the pseudocode overview of code functionality in Section 3.

- To execute the supplied code refer to Section 2 of this document.
- To replicate this code refer to Sections 3 to 13 of this document.

### 1.1 Hardware and Operating System

We utilized a MacBook Pro with an Apple M1 Pro chip, featuring an 8-core CPU and 16 GB of memory. The startup disk is the Macintosh HD. The system operates on macOS 13.2.1 (22D68) **Ventura**.

### 1.2 Software Versions and Dependencies

We used R version 4.3.1 (2023-06-16) **Beagle Scouts** for the x86\_64-apple-darwin20 platform. Furthermore, we used the RStudio environment, version 2023.06.2+561 to execute the our code. Table 1 below shows packages and versions utilized in this code.

Table 1: Package Version

Package	Version
Rgraphviz	2.44.0
bnlearn	4.8.3
gRain	1.3.13
visNetwork	2.1.2
ggplot2	3.4.2
zoo	1.8.12
scales	1.2.1
gridExtra	2.3
dplyr	1.1.2
MASS	7.3.60
svglite	2.1.1
tidyverse	2.0.0

### 1.3 Install Time

Installation of RStudio and R is contingent upon the specific hardware and operating system in use. This also holds true for the packages employed within the code. However, the packages used in this code typically take a few minutes to install.

### 1.4 Run Time

The run time on our operating system for this code is approximately 15 seconds.

## 2 Full code execution instructions

Following the instruction below, you can execute our code for the 5-year lag period and generate DAGs and conditional dependency graphs. This code is in turn broken up into an annotated and justified workflow in Sections 5 to 13. The expected outputs for this code are visualized DAGs presented in the **viewer** tab in RStudio and conditional dependency plots are saved as **.svg** files in the **repository** folder. Additional results are printed in the **console** in RStudio.

1. Download the **repository** files available in GitHub via this link..
2. Unzip the **repository** files.
3. Open the **Hydro EIA Code 5-year model.R** and **Graph.Generator.R** in RStudio.

4. Install any required packages and load them, refer to Section 5.1.
5. Set the working directory to the **repository** folder. See Section 5.2. Ensure the working directory in both R files is correctly set.
6. Execute the code.
7. The conditional dependency figures will be saved automatically in the **repository** folder. The DAGs will be displayed in the **viewer**, all other results are displayed in the **console**.

### 3 Pseudocode Overview

Here we provide the pseudocode overview of our algorithm.

---

**Algorithm 1** Pseudocode Overview of Code Functionality

---

```
1: START
2: INITIALIZE the R environment:
3:   a. Install and load required packages.
4:   b. Set the working directory.
5: DEFINE required functions:
6:   a. DAG visualizer.
7:   b. Box-Cox transformer.
8:   c. Goodness of fit test.
9:   d. D-separation test.
10: PREPARE the data:
11:   a. Generate delta and lags.
12:   b. Create averaged and lagged variables.
13:   c. Ensure data is numeric.
14:   d. Check for Gaussian distribution.
15:   e. Transform non-Gaussian variables.
16:   f. Discretize variables.
17:   g. Create the final data-frame.
18: CREATE a blacklist.
19: VISUALIZE the hypothesized DAG.
20: CONSTRUCT score-based DAGs.
21: PLOT conditional dependency results.
22: TEST goodness of fit for each node.
23: PERFORM d-separation analysis.
24: END
```

---

## 4 Algorithm Overview

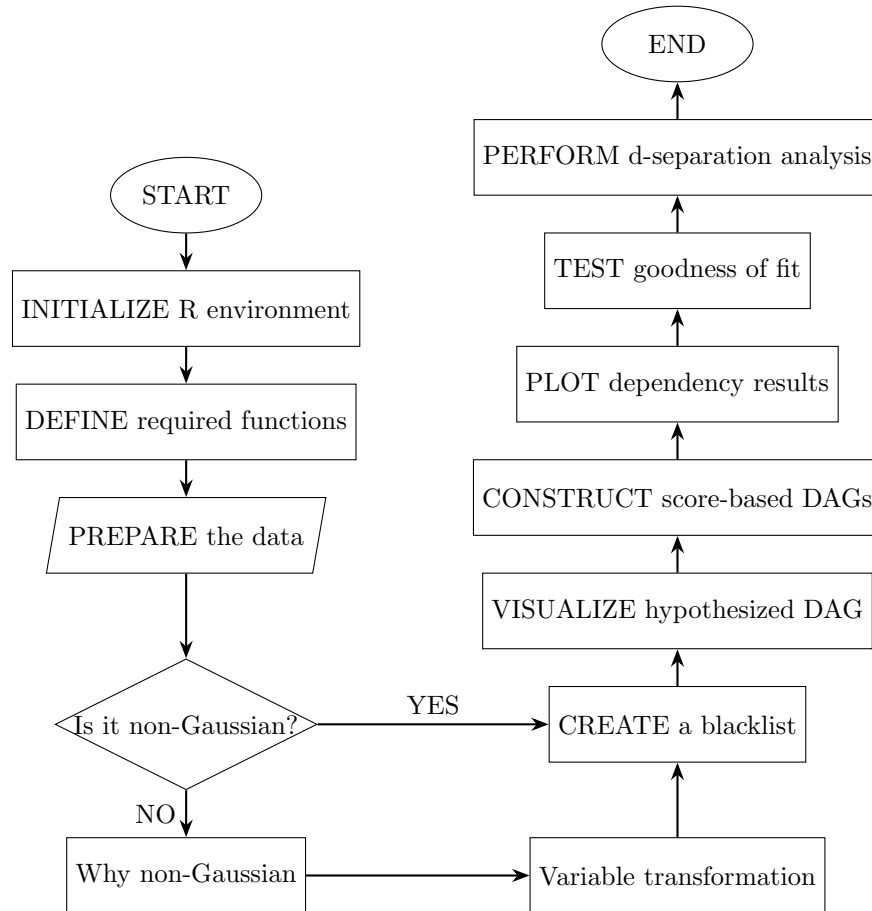


Figure 1: Algorithm to construct and evaluate BN models.

## 5 Preparing the R Environment

### 5.1 Install and Load the Required Packages

We use the following packages throughout the code. The core package that must be installed is the `bnlearn` package<sup>1</sup>.

```

1 # Install Rgraphviz from Bioconductor
2 if (!requireNamespace("BiocManager", quietly = TRUE))
3   install.packages("BiocManager")
4 BiocManager::install("Rgraphviz")
5
6 # List of other packages to install
7 packages <- c("bnlearn", "gRain", "visNetwork", "ggplot2",
8               "zoo", "scales", "gridExtra", "dplyr", "MASS", "svglite", "tidyverse")
9
10 # Install packages
11 install.packages(packages)
12
13 # Load all the required packages
14 invisible(lapply(c("Rgraphviz", "bnlearn", "gRain", "visNetwork", "ggplot2",

```

<sup>1</sup>Scutari, M., Silander, T., and Ness, R. (2023). Bayesian Network Structure Learning, Parameter Learning and Inference (4.8.3) [R]. Available: <https://www.bnlearn.com/>

```

15         "stats", "zoo", "scales", "gridExtra", "dplyr", "MASS", "svglite", ...
        "tidyverse"), library, character.only = TRUE))

```

## 5.2 Set Working Directory

Set the working directory and read in the data. Make sure to adjust the directory path to match your setup.

```

1 # Set Working Directory (change for your setup)
2 setwd("/Users/Documents/Hydro EIA Code")
3
4 # Read the data-set
5 hydro.data <- read.csv("hydro_var_aug23.csv")

```

## 6 Required Functions

We have created the following functions that will assist in creation and evaluation of our causal directed acyclic graphs (DAGs).

### 6.1 DAG Visualizer

Visualizes and returns DAGs created by the `bnlearn` package.

```

1 plot.network <- function(structure, ht = "400px", title){
2
3   # Unique nodes from the arcs of the structure are identified.
4   nodes.uniq <- unique(c(structure$arcs[,1], structure$arcs[,2]))
5
6   # A data frame for nodes is created with attributes like id, label, color, and shadow.
7   nodes <- data.frame(id = nodes.uniq,
8                       label = nodes.uniq,
9                       color = "maroon",
10                      shadow = TRUE)
11
12   # A data frame for edges is created with attributes like source, target, arrow ...
13   # direction, and other visual properties.
14   edges <- data.frame(from = structure$arcs[,1],
15                      to = structure$arcs[,2],
16                      arrows = "to",
17                      smooth = TRUE,
18                      shadow = TRUE,
19                      color = "black")
20
21   # The network is visualized using the visNetwork function and returned.
22   return(visNetwork(nodes, edges, height = ht, width = "100%"))
23 }

```

### 6.2 Box-Cox Transformer

Evaluates the data-set and checks for normality (using Shapiro-Wilk test), transforms the non-Gaussian variables using Box-Cox transformation. Re-evaluates the transformed variables with the Shapiro-Wilk test and checks for normality. Returns the results.

```

1 transform_and_test <- function(df, non_gaussian_vars){
2
3   # Lists to store variables that remain non-Gaussian after transformation and those that ...
4   # were successfully transformed.
5   still_non_gaussian <- vector("list")
6   transformed_vars <- vector("list")
7   df_new = df
8
9   # Each non-Gaussian variable is processed.
10  for (var in non_gaussian_vars) {
11
12    # The minimum value of the variable is determined.
13    min_value <- min(df[[var]], na.rm = TRUE)
14
15    # If the minimum value is less than or equal to zero, a constant is added to make it ...
16    # positive.
17  }
18 }

```

```

15   if (min.value ≤ 0) {
16     constant <- abs(min.value) + 1
17     df[[var]] <- df[[var]] + constant
18   }
19
20   # The Box-Cox transformation parameter (lambda) is estimated.
21   bc <- boxcox(df[[var]] ~ 1, plotit = FALSE)
22   lambda <- bc$x[which.max(bc$y)]
23
24   # Depending on the value of lambda, the variable is transformed.
25   if(abs(lambda) ≤ 1e-5){
26     transformed_var <- log(df[[var]])
27   } else {
28     transformed_var <- (df[[var]]^lambda - 1) / lambda
29   }
30
31   # The transformed variable is tested for normality using the Shapiro-Wilk test.
32   shapiro.test <- shapiro.test(transformed_var)
33
34   # Based on the p-value, the variable is categorized as still non-Gaussian or ...
35   # successfully transformed.
36   if (shapiro.test$p.value < 0.05) {
37     still_non_gaussian <- c(still_non_gaussian, var)
38   } else {
39     df_new[[var]] <- transformed_var
40     transformed_vars <- c(transformed_vars, var)
41   }
42
43   # The modified data frame, list of still non-Gaussian variables, and list of transformed ...
44   # variables are returned.
45   list(df = df_new, still_non_gaussian = still_non_gaussian, transformed = transformed_vars)

```

### 6.3 Goodness of Fit Test

Evaluates each variable's goodness of fit for the DAGs produced by the **bnlearn** package.

Function to compute the  $R^2$  (r-squared) metric for continuous variables, representing the proportion of variance in the dependent variable that's predictable from the independent variable.

```

1 evaluate_fit_continuous <- function(actual, predicted) {
2
3   # An empty list for metrics is initialized.
4   metrics <- list()
5   # The r-squared metric is computed and stored in the metrics list.
6   metrics$rsquared <- 1 - sum((predicted - actual)^2) / sum((actual - mean(actual))^2)
7   # The metrics list is returned.
8   return(metrics)
9 }

```

Function to compute the accuracy metric for discrete variables, defined as the ratio of correctly predicted values to the total number of values.

```

1 evaluate_fit_discrete <- function(actual, predicted) {
2
3   # An empty list for metrics is initialized.
4   metrics <- list()
5   # The accuracy metric is computed and stored in the metrics list.
6   metrics$accuracy <- sum(actual == predicted) / length(actual)
7   # The metrics list is returned.
8   return(metrics)
9 }

```

### 6.4 D-Separation Test

Uses conditional dependence test to evaluate relationships between independent nodes in the DAG produced by the **bnlearn** package while factoring in the available data.

We have created the `dsep.dag` function that can evaluate d-separation for any given node pair. This function uses the optimized DAG results to identify the following for each node pair and then calculates conditional independence: 1. Parents, 2. Neighbors (i.e., parents and children for each node) and 3. Markov-Blanket (i.e., parents, children and parents of children for each node). Furthermore, this function uses data to evaluate d-separation in addition to the results of DAG discovery. This combines the functionalities of the `ci.test` and `dsep` functions available in the `bnlearn` package. Where, `ci.test` exclusively utilizes data, while `dsep` solely employs DAGs.

#### 6.4.1 Function Usage

The `dsep.dag` function is used as follows:

```
dsep.dag(x, data, z)
```

#### 6.4.2 Parameters

The `dsep.dag` function accepts the following parameters:

- **x**: an object of class `bn`
- **data**: a data frame containing the variables in the model
- **z**: a list, where each element is a character vector representing a pair of node labels

Available conditioning sets (that are printed automatically) are:

1. **parents**: set of causal variables for each node
2. **neighbors**: set of causal and response variables for each node
3. **markov\_blanket**: set of causal and response variables including causal variables of responses for each node

#### 6.4.3 Return Value

The `dsep.dag` function returns the following for each node pair:

- **Results for parents**: a character string
- **Results for neighbors**: a character string
- **Results for markov\_blanket**: a character string

The character string can yield one of the following values, representing the outcome of the d-separation analysis:

1. **Conditionally Independent**: The variables are independent given the observed variables.
2. **Potential Missing Link**: There might be a missing link between the variables.
3. **Uncertain - Further Analysis Needed**: The relationship between the variables is uncertain, necessitating further analysis.
4. **NA**: The node pair is either not available in the DAG, or they are connected, implying conditional dependence.



#### 6.4.4 Underpinning Functions

To construct and utilize the `dsep.dag` function, we implemented the following functions sequentially:

Function to check if a directed arc exists between two nodes in a given set of arcs. This eliminates manually checking the arcs in an existing DAG.

```
1 arc_exists <- function(from, to, existing_arcs) {
2   # The existence of the arc is checked and the result is returned.
3   return(any(existing_arcs[existing_arcs[, "from"] == from, "to"] == to))
4 }
```

Function to perform d-separation tests on all non-adjacent pairs of nodes in a given DAG using the dataset provided.

```
1 perform_dsep_tests <- function(dag, data) {
2   # Existing arcs in the DAG are extracted.
3   existing_arcs <- arcs(dag)
4
5   # An internal function to check for arc existence is defined.
6   arc_exists <- function(from, to, existing_arcs) {
7     return(any(existing_arcs[existing_arcs[, "from"] == from, "to"] == to))
8   }
9
10  # A list for non-adjacent node pairs is initialized.
11  non_adj_pairs <- list()
12
13  # Non-adjacent node pairs are identified.
14  for (node1 in node.ordering(dag)) {
15    for (node2 in node.ordering(dag)) {
16      if (!arc_exists(node1, node2, existing_arcs) && !arc_exists(node2, node1, ...
17        existing_arcs) && node1 != node2) {
18        non_adj_pairs <- append(non_adj_pairs, list(c(node1, node2)))
19      }
20    }
21  }
22
23  # A list for test results is initialized.
24  results_parents <- list()
25  results_nbr <- list()
26  results_mb <- list()
27
28  # D-separation tests are performed for each non-adjacent pair.
29  for (pair in non_adj_pairs) {
30    # Conditioning on parents of a node
31    conditioning_set <- setdiff(unique(c(bnlearn::parents(dag, pair[1]), ...
32      bnlearn::parents(dag, pair[2]))), pair) # Combined Markov blanket excluding x and y
33    test <- ci.test(pair[1], pair[2], conditioning_set, data = data)
34    results_parents[[paste0(pair[1], "-", pair[2])]] <- test$p.value
35    # Conditioning on immediate neighbors of a node
36    conditioning_set <- setdiff(unique(c(nbr(dag, pair[1]), nbr(dag, pair[2]))), pair) # ...
37      Combined Markov blanket excluding x and y
38    test <- ci.test(pair[1], pair[2], conditioning_set, data = data)
39    results_nbr[[paste0(pair[1], "-", pair[2])]] <- test$p.value
40    # Conditioning on Markov Blanket of a node
41    conditioning_set <- setdiff(unique(c(mb(dag, pair[1]), mb(dag, pair[2]))), pair) # ...
42      Combined Markov blanket excluding x and y
43    test <- ci.test(pair[1], pair[2], conditioning_set, data = data)
44    results_mb[[paste0(pair[1], "-", pair[2])]] <- test$p.value
45  }
46
47  return(list(parents = results_parents, neighbors = results_nbr, markov_blanket = ...
48    results_mb))
49 }
```

Function to categorize the p-values from d-separation tests into three categories: "Conditionally Independent", "Potential Missing Link", and "Uncertain - Further Analysis Needed".

```

1 interpret_dsep_pvalues <- function(pvalues, threshold_low = 0.05, threshold_high = 0.95) {
2   # P-values are categorized based on predefined thresholds.
3   categories <- sapply(pvalues, function(p) {
4     if (p > threshold_high) {
5       return("Conditionally Independent")
6     } else if (p < threshold_low) {
7       return("Potential Missing Link")
8     } else {
9       return("Uncertain - Further Analysis Needed")
10    }
11  })
12  return(categories)
13 }

```

Function to organize and print the interpret\_dsep\_pvalues function results.

```

1 interpret_print <- function(node1, node2, interpretations_list) {
2   # Constructing the key string from the node names
3   key <- paste0(node1, "-", node2)
4
5   interpretations <- lapply(interpretations_list, function(interpretation) {
6     interpretation[key]
7   })
8   names(interpretations) <- c("parents", "neighbors", "markov_blanket")
9   # Printing the results
10  cat("From", node1, "To", node2, ":\n")
11  for (name in names(interpretations)) {
12    cat("Results for", name, ":\n")
13    cat(interpretations[[name]], "\n\n") # Directly prints the value without the key
14  }
15 }

```

Function that wraps all of the previous functions into one concise function, i.e., the dsep.dag function.

```

1 dsep.dag <- function(dag, data, node_pairs) {
2   # Performing D-separation tests
3   dsep_results <- perform_dsep_tests(dag, data)
4
5   # Translating the results
6   interpretations_list <- lapply(dsep_results, interpret_dsep_pvalues)
7
8   # Initializing a list to hold the interpreted results for each node pair
9   interpreted_results <- list()
10
11  for (node_pair in node_pairs) {
12    key <- paste0(node_pair[1], "-", node_pair[2])
13    interpret_print(node_pair[1], node_pair[2], interpretations_list)
14    interpreted_results[[key]] <- lapply(interpretations_list, function(interpretation) ...
15      interpretation[key])
16  }
17
18  # Returning the list of interpreted results
19  return(interpreted_results)
20 }

```

## 7 Data Preparation

Modify the variables as needed, we introduced lags and deltas to ensure we capture the temporal difference of causality for our variables. Refer to the main document for further information. Note that in this section we only show the code for some variables transformations. Refer to the main document for further details for each transformation.

### 7.1 Generating delta and lags

We manually introduced deltas and lags into our variables based on the hypothesized model.

```

1 {
2   # New intertie capacity in every year
3   hydro.data$INTERTIE_new = NA
4   for(i in 2:nrow(hydro.data)){
5     hydro.data$INTERTIE_new[i] = hydro.data$INTERTIE[i] -
6       hydro.data$INTERTIE[i-1]
7   }
8   # Sum of new intertie capacity in preceding 5 years
9   hydro.data$INTERTIE_5y = NA
10  for(i in 5:nrow(hydro.data)){
11    hydro.data$INTERTIE_5y[i] = sum(hydro.data$INTERTIE_new[(i-4):i], na.rm=T)
12  }
13  # Sum of new intertie capacity in preceding 5 years lagged by 5 years
14  hydro.data$INTERTIE_5y_lag_5y = NA
15  for(i in 10:nrow(hydro.data)){
16    hydro.data$INTERTIE_5y_lag_5y[i] = hydro.data$INTERTIE_5y[i-5]
17  }
18 }

```

## 7.2 Creating averaged and lagged variables

We used loops to generate new variables that represent 5-year averages and their respective lags.

```

1 lag_periods <- c(5)
2 new_vars <- c("INTERTIE", "INSTALLED", "DEMAND_QC", "DEMAND_US", "INVESTMENT", "EXPORTS", ...
3   "PRICE")
4 for (var in new_vars) {
5   var_new <- paste0(var, "_new")
6
7   for (lag_period in lag_periods) {
8     # moving average
9     var_avg <- paste0(var_new, "_avg_", lag_period, "y")
10    hydro.data[[var_avg]] <- zoo::rollapplyr(hydro.data[[var_new]], width = lag_period, ...
11      FUN = mean, fill = NA)
12
13    # lagged average
14    var_lag_avg <- paste0(var_avg, "_lag_", lag_period, "y")
15    hydro.data[[var_lag_avg]] <- dplyr::lag(hydro.data[[var_avg]], lag_period)
16  }
17 }

```

## 7.3 Creating a subset for averaged and lagged variables

We subset our data-set to only retain 5-year averaged and lagged variables.

```

1 vars.exclude.5y = c(1:8, grep("_new$", colnames(hydro.data)))
2 hydro.data.subset.5y =
3   hydro.data[, setdiff(1:ncol(hydro.data),
4     vars.exclude.5y)]

```

## 7.4 Ensuring data-set is numeric

It is critical to ensure that the data-set is numeric before proceeding to perform any further analysis.

```

1 for(i in 1:ncol(hydro.data.subset.5y)){
2   hydro.data.subset.5y[,i] = as.numeric(hydro.data.subset.5y[,i])
3 }

```

## 7.5 Creating a new data-frame

After preparing our data-set we create a new data-frame that excludes rows with missing values. It must be noted that introducing lags and deltas will result in some information loss, therefore selecting an optimum lag period is imperative.

```

1 # Create new dataframe
2 df.5y = hydro.data.subset.5y

```

```

3
4 # Creating new data frames with minimum number of rows cut off
5 lag.cols.5y = grep("lag", colnames(df.5y))
6 df.5y.no.lags = df.5y[, setdiff(1:ncol(df.5y),
7                               lag.cols.5y)]
8 df.5y.rows.to.cut = which(apply(df.5y, 1, function(x) sum(is.na(x)) > 0))
9 df.5y.no.lags.rows.to.cut = which(apply(df.5y.no.lags, 1, function(x) sum(is.na(x)) > 0))
10 df.5y.with.lags.no.NA = df.5y[setdiff(1:nrow(df.5y),
11                                     df.5y.rows.to.cut), ]

```

## 7.6 Check for Gaussian distribution

We perform the Shapiro-Wilk test to check for normality across our variables. Bayesian models we use in this analysis require normality to ensure an accurate estimation. Variables that don't pass the Shapiro-Wilk test a p-value less than the specified significance level (0.05) are considered non-Gaussian and are stored in a list.

```

1 significance.level <- 0.05
2 non_gaussian_5y_lag <- vector("list")
3 for (var in colnames(df.5y.with.lags.no.NA)) {
4   # Shapiro-Wilk Test
5   shapiro.test <- shapiro.test(df.5y.with.lags.no.NA[[var]])
6   print(paste("Shapiro-Wilk Test for", var, "- p-value:", shapiro.test$p.value))
7
8   if (shapiro.test$p.value < significance.level) {
9     non_gaussian_5y_lag <- c(non_gaussian_5y_lag, var)
10  }
11 }

```

## 7.7 Transforming non-Gaussian variables

We transform the variables that are non-Gaussian, obtained from the previous step. We use the 'transform\_and\_test' function created in Section 6.2 to transform the variables using Box-Cox transformation. If the variables are still non-Gaussian after transformation the function reverts them to their original state.

```

1 result_5y_lag <- transform_and_test(df.5y.with.lags.no.NA, non_gaussian_5y_lag)
2 df.5y.with.lags.no.NA <- result_5y_lag$df
3 still_non_gaussian_5y_lag <- result_5y_lag$still_non_gaussian

```

## 7.8 Discretizing variables

Variables that cannot be transformed using the Box-Cox transformation are then discretized into categories based on their values and histograms.

```

1 df.5y.with.lags.no.NA$EXPORTS_new_avg_5y <- cut(df.5y.with.lags.no.NA$EXPORTS_new_avg_5y, ...
  breaks = c(min(df.5y.with.lags.no.NA$EXPORTS_new_avg_5y, na.rm=T), 0, ...
  max(df.5y.with.lags.no.NA$EXPORTS_new_avg_5y, na.rm=T)), labels = c("negative", ...
  "positive"), include.lowest = TRUE, ordered.result = TRUE)

```

## 7.9 Creating the final data-frame

We drop the variables that are identical after transformations. Finally we create the final data-frame by subsetting the current data-frame.

```

1 # Drop the variables
2 var.drop = c("INSTALLED_new_avg_5y", "INTERTIE_new_avg_5y", "INSTALLED_new_avg_5y_lag_5y", ...
3             "INTERTIE_new_avg_5y_lag_5y", "INSTALLED_lag_5y", "INTERTIE_lag_5y")
4 # Subset the dataset
5 df.5y.with.lags.no.NA =
6   df.5y.with.lags.no.NA[, setdiff(1:ncol(df.5y.with.lags.no.NA),
7                                   match(var.drop, colnames(df.5y.with.lags.no.NA)))]
8
9 # Selecting columns of interest
10 selected.columns <- c('INSTALLED_5y_lag_5y', 'INSTALLED_5y', ...
11                      'DEMAND_QC_new_avg_5y_lag_5y', 'INVESTMENT_5y', 'PRICE_5y_lag_5y',
12                      'INTERTIE_5y_lag_5y', 'INTERTIE_5y', 'DEMAND_US_new_avg_5y_lag_5y', ...
13                      'PRICE_5y',

```

```

11         'EXPORTS-new-avg-5y-lag-5y', 'EXPORTS-5y', 'DEMAND-QC-5y', ...
12         'DEMAND-US-5y')
13 # Subsetting the dataframe and creating the final dataframe
14 df.expert.5y <- df.5y.with.lags.no.NA[, selected.columns]

```

## 8 Creating a Blacklist

We have to set the parameters and boundaries for our Bayesian model. We do this by creating an allowable list of arcs derived from our hypothesized "expert" model. Using a for-loop we then convert this allowable list to a blacklist.

```

1 # The allow list is initialized with specific variable pairs.
2 allow.list.expert =
3   data.frame(matrix(c(
4     # Various variable pairs are listed here.
5     "INVESTMENT-5y", "INSTALLED-5y",
6     ...
7     "DEMAND-US-5y", "PRICE-5y"),
8     ncol = 2, byrow=TRUE))
9
10 # Column names for the allow list are assigned.
11 colnames(allow.list.expert) = c("From", "To")
12
13 # The black list is initialized with a placeholder value.
14 black.list.expert = NA
15
16 # For each pair of variables in the final data-frame, a check is performed.
17 # If the pair is not found in the allow list, it is added to the black list.
18 for(i in 1:ncol(df.expert.5y)){
19   for(j in 1:ncol(df.expert.5y)){
20     from.test = colnames(df.expert.5y)[i]
21     to.test = colnames(df.expert.5y)[j]
22
23     if(length(which(allow.list.expert$From==from.test&
24                   allow.list.expert$To==to.test))==0){
25       black.list.expert =
26         rbind(black.list.expert, c(from.test, to.test))
27     }
28   }
29 }
30
31 # Column names for the black list are assigned.
32 colnames(black.list.expert) = c("From", "To")
33
34 # The placeholder value in the black list is removed.
35 black.list.expert = black.list.expert[2:nrow(black.list.expert),]

```

## 9 Visualizing the Hypothesized DAG

We construct the hypothesized "expert" DAG using the 'model2network' function from the **bnlearn** package. We visualize this DAG using 'plot.network' function that we created in Section 6.1. Figure 2 shows the visualized expert DAG.

```

1 # A DAG is constructed using expert knowledge.
2 dag.expert.5y <- model2network(" [INSTALLED-5y.lag-5y] [DEMAND-QC-new-avg-5y.lag-5y]
3 [PRICE-5y.lag-5y] [INTERTIE-5y.lag-5y] [DEMAND-US-new-avg-5y.lag-5y] [EXPORTS-new-avg-5y.lag-5y]
4 [DEMAND-QC-5y] [DEMAND-US-5y]
5 [INSTALLED-5y|DEMAND-QC-new-avg-5y.lag-5y:INVESTMENT-5y:PRICE-5y.lag-5y:INTERTIE-5y.lag-5y]
6 [INTERTIE-5y|INSTALLED-5y.lag-5y:INVESTMENT-5y:DEMAND-US-new-avg-5y.lag-5y:PRICE-5y.lag-5y]
7 [INVESTMENT-5y|DEMAND-QC-new-avg-5y.lag-5y:EXPORTS-new-avg-5y.lag-5y]
8 [EXPORTS-5y|INTERTIE-5y:INSTALLED-5y:PRICE-5y]
9 [PRICE-5y|DEMAND-QC-5y:DEMAND-US-5y]")
10
11 # The constructed DAG is visualized with a specified height.
12 plot.network(dag.expert.5y, ht = "600px")

```

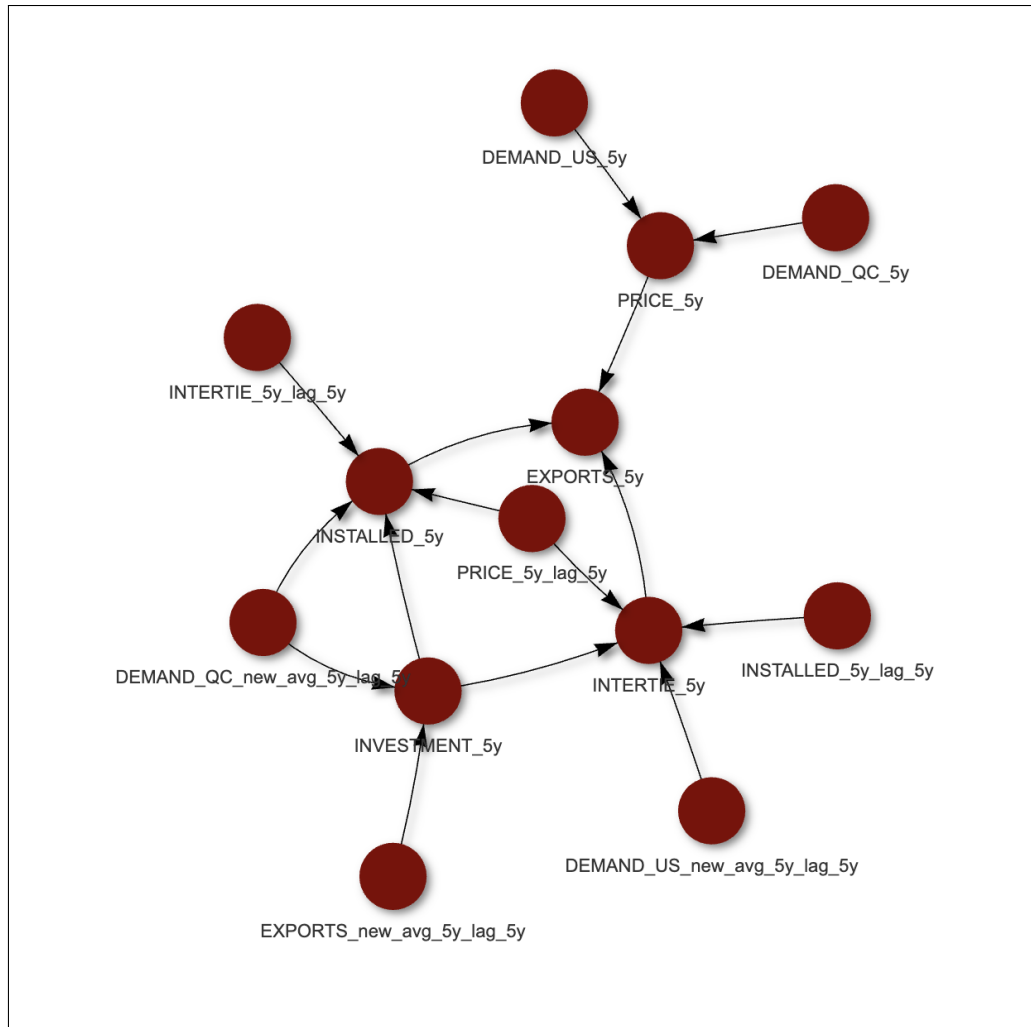


Figure 2: Hypothesized DAG visualization.

## 10 Constructing Score-Based DAGs

We construct DAGs using the score-based hill-climb (HC) algorithm. We use the 'hc' and 'bn.fit' functions from the `bnlearn` package. We visualize this DAG using 'plot.network' and 'graphviz.chart' functions. We use three different scoring functions; 1. 'loglik-cg', 2. 'aic-cg' and 3. 'bic-cg'.

```

1 # DAG created using the loglik-cg score function and HC algorithm
2 dag.expert.5y.emp <- hc(df.expert.5y, score = "loglik-cg", blacklist = black.list.expert, ...
   debug = FALSE)
3 par(mar=c(1,1,1,1))
4 # Fitting the model
5 model.expert.5y.emp = bn.fit(dag.expert.5y.emp, df.expert.5y)
6 #Visualizing model's conditional probabilities using the graphviz.chart
7 graphviz.chart(model.expert.5y.emp, type = "barprob", grid = TRUE, bar.col = "darkgreen",
   strip.bg = "lightskyblue")
8
9 dev.off()
10 # Network visualized using plot.network
11 plot.network(dag.expert.5y.emp, ht = "600px")
12
13 # DAG created using the aic-cg score function and HC algorithm
14 dag.expert.5y.emp.aic <- hc(df.expert.5y, score = "aic-cg", blacklist = black.list.expert)
15 plot.network(dag.expert.5y.emp.aic, ht = "600px")
16 par(mar=c(1,1,1,1))

```

```

17 model.expert.5y.emp.aic = bn.fit(dag.expert.5y.emp.aic, df.expert.5y)
18 graphviz.chart(model.expert.5y.emp.aic, type = "barprob", grid = TRUE, bar.col = "darkgreen",
19               strip.bg = "lightskyblue")
20 dev.off()
21
22 # DAG created using the bic-cg score function and HC algorithm
23 dag.expert.5y.emp.bic <- hc(df.expert.5y, score = "bic-cg", blacklist = black.list.expert)
24 plot.network(dag.expert.5y.emp.bic, ht = "600px")
25 par(mar=c(1,1,1,1))
26 model.expert.5y.emp.bic = bn.fit(dag.expert.5y.emp.bic, df.expert.5y)
27 graphviz.chart(model.expert.5y.emp.bic, type = "barprob", grid = TRUE, bar.col = "darkgreen",
28               strip.bg = "lightskyblue")
29 dev.off()

```

## 11 Plotting Conditional Dependency Results

We used the conditional probability results from the DAG model to evaluate the relationship between various nodes. We used visualization to plot these dependencies. This was achieved by creating a separate R code called `Graph_Generator.R` located in the same directory as the main code. This code is then recalled using the `source` function.

It must be noted that the type of each plot is dependant on the type of each node and the number of its parents. Therefore, here we only present one node's results as an example. Refer to the supplemental information document for further details for each node. Figure 3 shows the results for this plot.

### 11.1 Running the Graph\_Generator.R Code

We use the following code to recall the `Graph_Generator.R` code from the directory. This will run that code in the background.

```

1 source("Graph_Generator.R")

```

### 11.2 Snippet of Graph\_Generator.R Code

Here is an snippet of `Graph_Generator.R` code where graphs are defined for nodes of interest and they're saved in the `.svg` format automatically when this code is called.

```

1 # 5. Conditional density: EXPORTS_5y | INSTALLED_5y + INTERTIE_5y + PRICE_5y
2 # Makind predictions for the node EXPORTS_5y using 'predict' function from the 'bnlearn' ...
   package.
3 df.expert.5y$EXPORTS_5y_pred = predict(model.expert.5y.emp, node = "EXPORTS_5y", data = ...
   df.expert.5y, method = "bayes-lw")
4 p<- ggplot(df.expert.5y, aes(x = PRICE_5y, y = INSTALLED_5y, size = EXPORTS_5y_pred, color ...
   = as.factor(INTERTIE_5y))) +
5   geom.point() +
6   labs(x = "PRICE",
7        y = "INSTALLED",
8        size = "Predicted EXPORTS",
9        color = "INTERTIE") +
10  scale_color_manual(values = palette)
11 ggsave(filename = "fig-8.svg", plot = p, device = "svg")

```

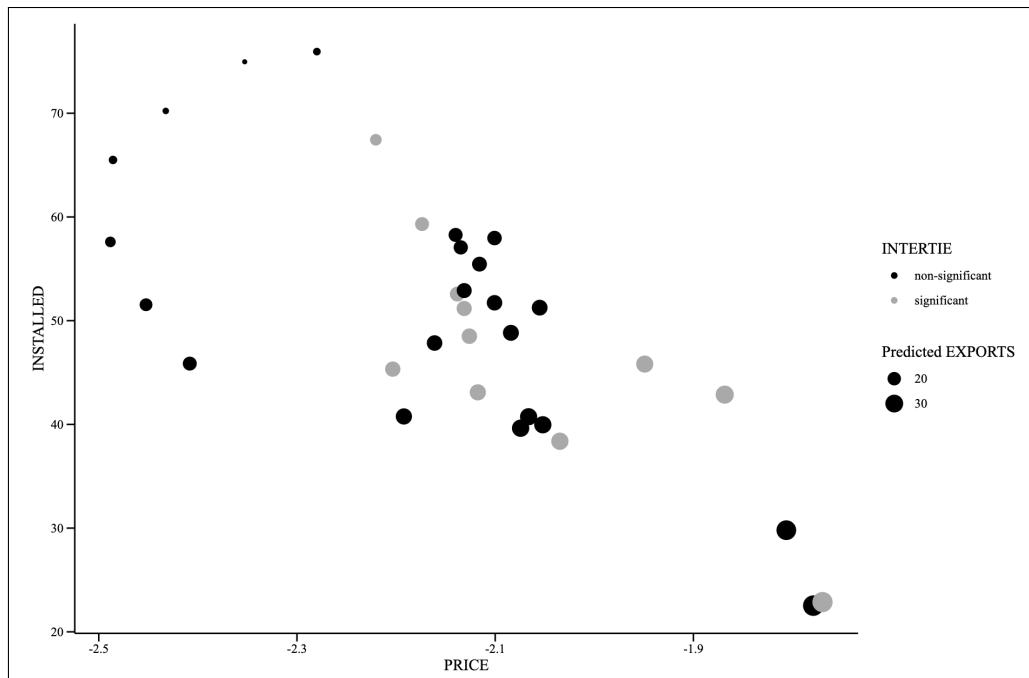


Figure 3: 5-year model's conditional dependency results for the exports node.

## 12 Goodness of Fit for Each Node

Here we use the `'evaluate_fit_discrete'` and `'evaluate_fit_continuous'` functions created in Section 6.3 to evaluate the goodness of fit for nodes with parent(s). We use the `predict()` function from the `bnlearn` package. Additionally we utilize Monte Carlo posterior inference method using the `bayes-lw` with 5000 parameters.

```

1 # identify discrete and continuous variables
2 discrete_vars <- c("INTERTIE.5y", "INVESTMENT.5y")
3 continuous_vars <- setdiff(colnames(df.expert.5y), c(discrete_vars, value = TRUE))
4 # Create a blank results list
5 results_loglik <- list()
6
7 # Looping over each node to generate its prediction using the 'predict' function from the ...
8   'bnlearn' package.
9 for (var in colnames(df.expert.5y)) {
10   if (!grepl("_pred$", var)) {
11     pred_column <- paste(var, "pred", sep = ".")
12     df.expert.5y[[pred_column]] <- predict(model.expert.5y.emp, node = var, data = ...
13       df.expert.5y, method = "bayes-lw", n = 5000)
14     actual_values <- df.expert.5y[[var]]
15     predicted_values <- df.expert.5y[[pred_column]]
16
17     # Calculate the goodness of fit for different type of variables using the functions we ...
18     # created previously
19     if (var %in% continuous_vars) {
20       results_loglik[[var]] <- evaluate_fit_continuous(actual_values, predicted_values)
21     } else if (var %in% discrete_vars) {
22       results_loglik[[var]] <- evaluate_fit_discrete(actual_values, predicted_values)
23     }
24   }
25 }

```

## 13 D-Separation Analysis

Here we use the `'arc_exists'`, `'perform_dsep_tests'` and `'interpret_dsep_pvalues'` functions created in Section 6.4 to evaluate the conditional dependence of nodes of interest if needed. This step helps resolve



any discrepancies or nuances that we might witness in different models.

```
1 # Defining the node_pairs list to identify nodes of interest where we want to perform ...  
  d-separation. Note that this list contains node in a "from", "to" format.  
2 different_edges <- compare(dag.expert.5y.emp, dag.expert.5y, arcs = TRUE)  
3 node_pairs <- different_edges$fp  
4 print(node_pairs)  
5 node_pairs <- lapply(seq_len(nrow(node_pairs)), function(i) as.character(node_pairs[i, ]))  
6  
7 # Using the dsep.dag function to calculate conditional dependency for each pair  
8 dsep_log <- dsep.dag(dag.expert.5y.emp, df.expert.5y, node_pairs)
```

