The DM Algorithm:

A Causal Search Algorithm for the Discovery of

MIMIC Models, with an Attempt to Recover a

Protein Signalling Network from a

High-Dimensional Ovarian Cancer Dataset

Alexander Murray-Watters

August 21, 2014

**Abstract**

Latent variables have long confounded attempts to determine causal structure when experiments cannot be conducted. While some methods exist for dealing with exogenous latent variables, endogenous latents remain neglected. This thesis presents a new algorithm (the DM algorithm) designed to discover causal structure for a restricted class of models when endogenous latents are present. The algorithm is non-parametric, and in simulations outperformed one of the most popular methods for handling endogenous latents (namely, factor analysis). As the DM algorithm is also capable of handling a surprising number of variables, the algorithm was run on a high-dimensional genomic dataset. Popular methods in genomics lack the ability to address large numbers of variables and provide less information about the latent structure than the DM algorithm, so this represents an improvement on the state-of-the-art.

# Acknowledgements

I'd like to begin by thanking my parents and twin for their tolerance of all my annoying questions over the years, as well as their willingness to listen while I worked through problems out loud.

Many thanks to Clark Glymour for allowing me to research a problem on which little progress might be made and for preventing me from falling down the various crevices that appear during research.

Many thanks also to Richard Scheines for his time and very helpful suggestions(such as on factor analysis), despite his new administrative burdens.

Needless to say, I'm also very grateful to both Clark and Richard for their willingness to be my advisors.

I'm grateful to Xinghua Lu for both providing the ovarian cancer dataset used in the data analysis and for kindly taking the time to answer questions about its contents.

I'd also like to thank Joe Ramsey for providing a great deal of assistance in dealing with some computational difficulties.

Finally, I'd like to express apprection of my fellow graduate students for all of the interesting and informative discussions we've had over the past year, as well as their support during stressful periods.

# Contents

# List of Figures

# List of Algorithms

# List of Source Code Files

# Chapter 1

# Introduction

One of the enduring features of science is the discovery of novel variables and their causal relations with each other and with known variables. The histories of physics, chemistry and biology all illustrate such discoveries, from novel fundamental particles, to atoms and their weights, to genes and the processes through which they produce proteins. As these examples illustrate, the novel variables are often unmeasured and unmeasurable at the time of their discovery. Similar issues of discovery arise in psychology, economics, and the many social sciences where the means of discovery, such as they are, have been primarily statistical, as in psychometrics.

In many of these domains an important research question concerns the identification of unmeasured variables that are causal intermediaries between measured variables, the selection of clusters of input and output variables that share a common intermediate, and the causal relations among the unmeasured intermediates. Causal hypotheses of these forms are often referred to as "MIMIC" models, short for Multiple Indicators Multiple Input Causes.

Several statistical techniques have been proposed or applied for the purpose

Figure 1.1: The structure of a MIMIC Model. $G_1 - G_4$ and $GE_1 - GE_4$ are observed variables, while $L_1 - L_2$ are unobserved variables.

of finding MIMIC models, each with its special limitations. The aim of this thesis is to describe, illustrate and prove correct (under restrictive assumptions) a new, non-parametric method for identifying a sub-class of MIMIC models. The thesis introduces prior methods used to discover hypothesized MIMIC models in data and explores their limitations. It then introduces the necessary conceptual and theoretical ideas, before discussing the new algorithm. This discussion includes both a worked example of the algorithm in practice, as well as a proof of correctness. Simulations are then discussed, comparing the performance of the new algorithm with that of factor analysis. A data analysis using the algorithm is then conducted on a genomic dataset and the results summarized. Finally, future avenues for research are discussed.

The very idea of algorithmic discovery of theories, or "models" has been challenged in the philosophical and statistical literatures. Notably, Carl Hempel, among the most eminent philosophers of science of the 20th century, denied the possibility on the very grounds that no algorithm could correctly discover

unmeasured, novel properties(Hempel, 1985). Developments in the late 20th century and this century show that Hempel's opinion should be clustered with Kant's endorsement of phlogiston in the first edition of his *Critique of Pure Reason*, and with Hegel's (possibly apocryphal) announcement that there are necessarily seven planets (Craig and Hoskin, 1992). Improvements in automated causal inference now appear almost monthly, and all manner of problems (cyclic feedback structures, latent variables, latent variables and feedback, non-linear systems, and non Gaussian variables) that were once thought insuperable have been solved.

## 1.1 MIMIC Models in Practice

A hypothesized MIMIC model consists of three parts: a set of unobserved (latent) variables, a set of observed effects of the latents (or "outputs"), and a set of observed causes of the latents (or "inputs").

The ability to discover a MIMIC structure is useful in a number of situations. In genomics, some researchers are interested in discovering how the effect of a genetic mutation propagates through the protein signalling network, resulting in (possibly) different observed frequencies of various proteins. Genes are read and transcribed into mRNA, which is then translated into various proteins. Genes influence both which, and in what quantity, proteins are made. A genetic mutation (say, one that causes cancer) can lead to differences in the resulting levels of protein species, protein combination, and folding configurations. How the effects of mutation propagate through the network is of fundamental biological importance.

In neuroscience, researchers are interested in understanding how a signal (i.e., a stimulus), propagates through the brain's neural,network, and ultimately how these processes produce behavior. fMRI resolution is about two seconds

(Logothetis, 2008)[pg. 3], and since the signal propagates through the network at a significantly faster rate, it is currently not possible to directly observe the propagation of a signal through the brain. An important problem is to reveal the network structure from indirect imaging data clustered into "Regions of Interest," and MIMIC models, if discoverable, offer the possibility of identifying unmeasured intermediates between Regions of Interest.

In order to estimate the size of the shadow economy (the portion of the economy not captured by GDP or other government statistics), economists have made use of MIMIC models. Bühn and Schneider (2008) used MIMIC models to examine economic loss attributed to the shadow economy in France. Giles (1999) employed a MIMIC model to create a time-series view of the shadow economy in New Zealand. Tedds (1998) also estimated a MIMIC model in order to determine the size of the shadow economy in Canada. DellAnno and Schneider (2006) published an article advocating further use of MIMIC models in economics.

In Lester (2008), the researcher used a MIMIC model to determine what factors related to the successful settlement of immigrants to Australia. Subjects were non-labor force participants, as well as economic and non-economic immigrants. Indicators of successful settlement included mental health, belief that the decision to migrate was correct, encouraging others to migrate to Australia, as well as reported level of life satisfaction.

Ríos-Bedoya et al. (2009) used MIMIC models to examine the strength of association between two latent factors (pleasant/unpleasant early smoking experiences) and current smoking status.

## 1.2 Current Methods for Discovering MIMIC Models

The most common method of specifying a MIMIC model is simply to make one up and test it statistically. As is well-known (Mayo, 1996), there are typically a multitude of alternative models that can pass standard tests of fit for a data set. We have no reason for confidence that the tacit search, whatever it is, that an investigator goes through in proposing a model is a reliable procedure that adequately considers the alternative. For that kind of confidence, statistical search methods are needed that can be shown to be at least asymptotically correct under explicit assumptions or, failing that, will at least search a broad space of alternative models.

Previously developed methods used to discover a posited MIMIC structure are problematic in several ways. In the case of factor analysis, the method is unreliable (which is clearly shown in the simulations reported in chapter 6). In others, the method is computationally intractable for most interesting problems (Markowetz et al., 2007; Tresch and Markowetz, 2008), or makes highly restrictive assumptions(Brodie, 2014).

The methods proposed to find MIMIC models range from the recent (nested effect models), to the cutting-edge (sparse endogenous latent search), to adaptations of the procedure that inaugurated algorithmic search early in the last century (factor analysis).

### 1.2.1 Factor Analysis

Factor Analysis models data ($\mathbf{P}$) by multiplying a principal component matrix ($\mathbf{w}$) by the projections of the data onto the principle components ($\mathbf{F}$). Any residual differences between the data and the result of the matrix multiplication

are accounted for by an error term, $\epsilon$. $\mathbf{X} = \mathbf{Fw} + \epsilon$

For example, imagine we have run factor analysis on a dataset (made up), and requested a model with 2 factors. We get back the following loading matrix:

|      | Factor1 | Factor2 |
|------|---------|---------|
| Var1 | 0.49    | 0.44    |
| Var2 | 0.16    | 0.82    |
| Var3 | 0.30    | 0.86    |
| Var4 | 0.21    | 0.36    |
| Var5 | 0.7     | 0.38    |
| Var6 | 0.78    | 0.12    |

A "loading" is the correlation between a variable and a given factor. So Var1 and Factor1 have a correlation of .49. The matrix is then converted to a graph by choosing some cutoff loading (commonly .3), and drawing an undirected edge between a factor (latent) and a variable only if their loading is greater than the cutoff. So in the example, Factor1 has edges to Var1, Var3, Var5,and Var6, while Factor2 has edges to Var1, Var2, Var3, Var4, and Var5. Figure 1.2 depicts an example resulting from such a conversion.

Factor analysis is often used to perform dimension reduction, or in other cases, to infer causal structures when there are unobserved (latent) variables present.

When being used to infer a causal structure, proponents of factor analysis distinguish between two kinds of analysis: exploratory and confirmatory. In exploratory factor analysis, the number of latents to be used is not decided prior to looking at the data. Instead, the data is used to select the number of latents (often by looking at a Scree plot). Confirmatory factor analysis begins with a number of different models already specified (i.e., the number of latents is already specified for each model), and uses a chi-square goodness of fit statistic to select which of the models will be the final model. This distinction between

Figure 1.2: An example of a MIMIC model found using factor analysis, from Hughes et al, 2010. The model depicts two latents (related to executive function). The outputs are three different test scales acquired when subjects were ages 4 and 6.

types of factor analysis based on prior specification is odd; it suggests that if a researcher were to test enough models using confirmatory analysis, the only difference in justification between the confirmatory and the exploratory version would be that the researcher wrote a list of models earlier. Yet practitioners are using completely different model selection criteria, even though there is little substantive difference between the two "kinds" of analysis. If we accept the distinction between types of factor analysis, however, then the factor analysis being discussed in this thesis is exploratory, not confirmatory.

The history of factor analysis (and its use to infer a causal structure) began when Charles Spearman observed that a collection of variables, specifically children's grades in different subjects, had a correlation matrix which followed a pattern of constraints (known as tetrad constraints). Using this pattern, he claimed that there exists a common latent "factor," which he referred to as G, representing a person's intelligence. Spearman's pattern failed to hold in general, which led to modifications by his students. Spearman's method was intractable at the time (Glymour et al., 1987). Thurstone later modified the method so that it was both computationally tractable and capable of handling more than one factor (Thurstone, 1934).

Unfortunately, factor analysis is an unreliable tool for causal inference, in that the method cannot reliably cluster variables around latents. This is due to the tendency of factor analysis to report a different structure when F is right multiplied by an orthogonal matrix transpose(m), and w is left multiplied by the same orthogonal matrix (m). This distressing result is known as the rotation problem (Shalizi, 2012). Simulations have been run which illustrate just how unreliable factor analysis is (see chapter 6).

15

Figure 1.3: Several examples of NEM models, from Tresch and Markowetz (2008). Note that the black dots are observed effects. X, Y, and Z are observed inputs. The black rectangles in each 3 by 12 matrix denote input variable influence, with row corresponding to input variable and column denoting output variable. Also note that the underlying latent structure is not specified.

### 1.2.2 Nested Effect Models (NEM)

Nested effect models (e.g., Markowetz et al. (2007); Tresch and Markowetz (2008)) have been developed specifically for genomics. They rely on perturbing genes (i.e., performing an intervention on a gene), observing the resulting changes in gene expression level, comparing these changes to the unperturbed levels, and finally fitting a model on the data (generally fitting a Bayesian network via a Markov-chain Monte Carlo type algorithm). Nested effect models are used to cluster genes with phenotypes (effects) (Markowetz et al., 2007) and order gene clusters based on "subset relations between phenotypes" (ibid). This is not the same as a MIMIC model, as causal connections between the latent parts of the signalling network are not contained in the reported NEM, leaving the details of their construction unspecified.

Despite recent improvements (Tresch and Markowetz, 2008), the various methods based on nested effect models are still computationally intractable for high dimensional problems (which some parts of genomics face). Additionally, the fastest method for constructing a NEM (Tresch and Markowetz, 2008) is very sensitive to the addition or removal of an edge in the model, leading to a very unstable likelihood function.

### 1.2.3 SELS (Sparse Endogenous Latent Search)

Another method recently proposed for studying causally incomplete systems (i.e., where relevant variables are unobserved) is known as SELS, or Sparse Endogenous Latent Search. It begins by calculating the rank of a dataset's "sparse factorization." The extent to which the resulting matrix is rank deficient gives the number of latent variables. This information is used to orient edges in a partial ancestral graph, the details of which are unimportant for this discussion.

Unfortunately, SELS makes a number of strong assumptions, preventing its use in the MIMIC model case. Specifically, in the true graph (i.e., the actual underlying causal structure), there can be no direct causal relations between latents, the parents of latents, or the children of latents (Brodie, 2014)[pg. 22]. Additional parametric assumptions, namely that variables follow a (multivariate) Gaussian distribution and are linearly related to one another (Brodie, 2014)[pg. 22], further reduce the method's applicability.

### 1.2.4 Find One Factor Clusters (FOFC)

Silva et al. (2006) developed a method for finding subsets of measured variables that share a single common unmeasured cause, and for estimating the causal relations among latent causes of different clusters. The procedure has been improved by Kummerfeld et al. (2014).This procedure can be used to find the latent structure of MIMIC models, with several limitations, including:

1. Relations between measured and latent variables must be linear.

2. Which inputs affect which latents is not usually identified.

3. Some measured variables will be eliminated.

4. Retained measured variables cannot influence other measured variables.

5. Each latent must have at least three output variables as children, neither of which is the child of another latent variable.

   The method has the advantage, however, that the true structure need not be singly connected.

# Chapter 2

# Causal Graphs

## 2.1  Background

In a causal graph, a directed edge is interpreted causally. The origin node of the edge is a cause of the terminating node (also known as a child of the origin node). The origin node is also referred to as the ancestor of the terminating node.

There are three kinds of basic structures used in causal graphs: Chains, forks, and colliders. Chains (Figure 2.1) consist of two or more nodes, with each node (except the last) possessing a child. Forks (Figure 2.3) consist of at least three nodes, with a common node causing all of the other nodes. Finally, a collider (Figure 2.2) consists of at least two nodes with a common child.

An undirected path connecting two variables (A and B) is a a series of (one or more) transitions between nodes, beginning with A and ending with B. A directed path is similar, except it only allows transitions to travel in the direction of the arrows.

**Definition 2.1.1.** (Causal Markov Condition) A node in a directed acyclic

Figure 2.1: A chain.



Figure 2.2: A collider.



Figure 2.3: A fork.

Figure 2.4: An example of a MIMIC model.

causal graph is independent of its non-descendants, given its parents.

**Definition 2.1.2.** (d-separation) Two nodes (X and Y) in a directed acyclic causal graph are d-separated (and therefore independent) by a set S, where X and Y are not elements of S, if and only if all paths connecting X and Y are blocked by S. If two nodes are not d-separated, then they are said to be d-connected (and are therefore dependent).

A path is blocked if the path contains an intermediate node which is neither a collider nor a descendant of a collider on the path, and said node has been conditioned on. A path is also blocked if there is a collider on the path (which hasn't been conditioned on) and no descendant of the collider has been conditioned on.

Finally, nodes displayed in a box are observed. In contrast, nodes in an oval or circle are unobserved, and are often referred to as "latents" or "factors".

To make both MIMC models and causal graphing language more concrete for the reader, here is a hypothetical example of a MIMIC model (Figure 2.4).

In the model, $L_2$ and $L_2$ are latents, as their nodes are represented by circles. $G_1$, $G_2$, $G_3$, $G_4$, $GE_1$, $GE_2$, $GE_3$, and $GE_4$ are all observed, as their nodes are represented by boxes.

As $L_1$ is a collider, and it blocks the only path between $G_1$ and $G_2$, $G_1$ is independent of $G_2$. As $L_1$ is a fork for $GE_1$ and $GE_2$, $GE_1$ and $GE_2$ are dependent. Similar relations hold for $G_3$, $G_4$, $GE_3$, and $GE_4$. The most interesting relation, however, is that while $G_1$ and $G_2$ are dependent with $GE_1$, $GE_2$, $GE_3$, and $GE_4$, $G_3$ and $G_4$ are only dependent with $GE_3$ and $GE_4$. This is due to $L_2$ acting as a collider, d-separating $G_3$ and $G_4$ from $GE_1$ and $GE_2$. It is also interesting to note that the independence of $G_3$ (or $G_4$) and $GE_1$ (or $GE_2$) disappears when $GE_3$ (or $GE_4$) is conditioned on, as doing so deactivates the collider at $L_2$, leaving $G_3$ and $GE_1$ d-connected.

Two other pieces of background information also need to be covered, as they are used in the algorithm discussed in the next chapter.

### 2.1.1    Sober's Criterion

Sober's criterion (Sober, 1998) is an empirical method for choosing between models which only cover cause and effect vs. ones which include latent variables. Figure 2.5 depicts two competing models which are representative of this problem. The criterion asks a simple question: Are the effects independent of one another when their causes are conditioned on? If the answer is yes, then the causes are not related to the effects via a latent variable. Otherwise, there is a latent present.

### 2.1.2    The PC Algorithm

The PC algorithm is a method for finding (given some assumptions) the set of causal graphs consistent with the data. It takes a dataset as input, and returns a

Figure 2.5: From Sober (1998). It depicts a "black box" model choice problem: Should the intermediate variable "thirst" be included or excluded from the model? Note that the bottom model is an example of a MIMIC model.

pattern consistent with the observed independence relations (i.e, an equivalence class of models). The pattern may or may not depict a unique causal graph, as a pattern can include undirected edges. The presence of an undirected edge in the pattern means that it (the edge's) direction is indeterminate, and one direction or the other (but not both) may be the truth. 1 is its pseudocode.

**Data**: Takes a dataset as input
**Result**: Returns a pattern.
A.) Form the complete undirected graph $C$ on the vertex set $V$.
B.)
$n = 0$.
**repeat**

    **repeat**

        select an ordered pair of variables $X$ and $Y$ that are adjacent in $C$ such that $Adjacencies(C, X) \setminus Y$ has cardinality greater than or equal to n, and a subset $S$ of $Adjacencies(C, X) \setminus Y$ of cardinality $n$, and if $X$ and $Y$ are d-separated given $S$ delete edge $X - Y$ from $C$ and record $S$ in $Sepset(X, Y)$ and $Sepset(Y, X)$;
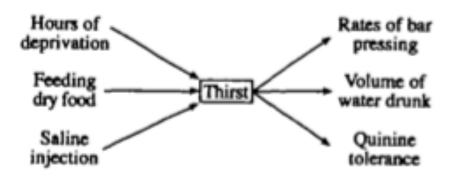
    **until** *All ordered pairs of adjacent variables $X$ and $Y$ such that $Adjacencies(C, X) \setminus Y$ has cardinality greater than or equal to n and all subsets $S$ of $Adjacencies(C, X) \setminus Y$ of cardinality n have been tested for d-separation*;

    $n = n + 1$;

**until** *For each ordered pair of adjacent vertices $X$, $Y$, $Adjacencies(C, X) \setminus Y$ is of cardinality less than n*;

C.) For each triple of vertices $X$, $Y$, $Z$ such that the pair $X$, $Y$ and the pair $Y$, $Z$ are each adjacent in $C$ but the pair $X$, $Z$ are not adjacent in $C$, orient $X - Y - Z$ as $X \rightarrow Y \leftarrow Z$ if and only if $Y$ is not in $Sepset(X, Z)$

D.) **repeat**

    If $A \rightarrow B$, $B$ and $C$ are adjacent, $A$ and $C$ are not adjacent, and there is no arrowhead at $B$, then orient $B - C$ as $B \rightarrow C$. If there is a directed path from $A$ to $B$, and an edge between $A$ and $B$, then orient $A - B$ as $A \rightarrow B$.

**until** *no more edges can be oriented*;

**Algorithm 1:** The pseudocode for the PC algorithm. (Spirtes et al., 2000)

The DM algorithm uses the PC algorithm for two tasks: classifying inputs and outputs using the edge direction found by PC, and determining which inputs are related to which outputs. In neither case is the output of PC used for its original purpose.

Having covered the basics of causal graphs, an example of a MIMIC model, and two other necessary pieces of background information, we can now move on to the algorithm itself.

# Chapter 3

# The DM (or detect.MIMIC) Algorithm

## 3.1   Definitions

Before delving into the algorithm, a number of definitions need to be made clear.

Inputs: Measured variables with only directed edges to (a) latent(s).

Outputs: Measured variables with a directed edge from a latent(s)

Singly connected: A graph is singly connected if there is at most a single undirected path connecting any pair of nodes.

Indegree: A variable's indegree is the number of directed edges pointing towards the variable,

Outdegree: A variable's outdegree is the number of directed edges pointing away from the variable.

Total Degree: The sum of a variable's indegree and outdegree.

Adjacency degree n: a function which returns the number of edges adjacent to node n.

Latent descendant: Given a latent (L1), A latent descendant is a latent variable caused by another latent variable.

## 3.2 Assumptions

Additionally, the algorithm currently requires a number of different assumptions.

Assumptions: The following properties are assumed to be true in the data generating process:

A1: Markov Assumption: Every variable is independent of its non-descendants given the variable's parents.

A2: Faithfulness: A graph and a probability distribution are faithful to one another if all the (un)conditional independence relations in the probability distribution are entailed by the graph and the Markov assumption.

A3: The true graph is acyclic.

A4: The true graph is singly connected.

A5: Every latent has at least two inputs and two outputs.

A6: No input has a path to an output except through a latent.

A7: Inputs are probabilistically independent of one another.
    Note: Generalizations of the algorithm are possible without this assumption (A7), but the information recovered may be reduced.

A8: Every measured variable is an input, an output, or a descendant of (an) output(s).

With those assumptions made, we can now move onto the algorithm itself.

It consists of seven steps, and is followed by a worked example.

## 3.3   Pseudocode

**Algorithm** : **DM(Data)**

**Data**: Takes a Dataset as Input
**Result**: Returns a MIMIC model
$PC :=$ A function returning the pattern produced by the PC algorithm.
**inputs** $:= NULL$ The set of inputs.
**outputs** $:= NULL$ The set of outputs.
$\mathbf{X} :=$ Data
$pc.pattern := PC(\mathbf{X}, depth = 0)$
$\mathbf{N} := Nodes(pcpattern)$
**for** *each n in* $\mathbf{N}$ **do**
    **if** $adjacency(n) \neq 0$ **then**
        **if** $adjacency(n) = outdegree(n)$ **then**
            | add $n$ to **inputs**
        **end**
        **else**
            | add $n$ to **outputs**
        **end**
    **end**
**end**
**Input**.**Parents**$(n) := PAR(n, pc.pattern) \cap$ **inputs**

**Algorithm 2:** Step 1 of the DM algorithm

**Latents** := $NULL$
**Latents**$(L)$ :=< **IN**$(L)$, **OUT**$(L)$, **LC**$(L)$ >
**for** *all L* **do**
  |   **Latents**$(L)$ :=< $NULL, NULL, NULL$ >
**end**
**Input**.**Parents** : The set of cluster assignments. Each member of
**Latents** (i.e., a specific latent) contains < **IN** = set of **inputs** for the
latent, **OUT** = set of **outputs**, and **LC** = set of latent children (i.e., a
latent descendant). > **for** *all x in* **outputs** **do**
    **if** *there exists a y in* **latents** *such that* **Input**.**Parents**$(x) = $ **IN**$(y)$
    **then**
    |   **OUT**$(y) := $ **OUT**$(y) \cup x$
    **end**
    **else**
        Create a new member, $z$, of **Latents**, with
        **Latents**$(z) :=< $ **IN**$(z) := $ **Input**.**Parents**$(x)$, **OUT**$(x) \cup x$,
        $NULL$ >
    **end**
**end**

**Algorithm 3:** Step 2 of the DM algorithm



**for** *each x, y in* **Latents** **do**
    **if** **IN**$(x)$ *is a proper subset of* **IN**$(y)$, *and* **IN**$(x)$ *is the largest such*
    *subset* **then**
        **LC**$(x) := $ **LC**$(x) \cup y$;
        **for** *all z in* **Latents** **do**
        |   **IN**$(z) := $ **IN**$(z) \setminus $ **IN**$(x)$
        **end**
    **end**
**end**

**Algorithm 4:** Step 3 of the DM algorithm



**for** *each x, y in* **Latents** **do**
    **if** **LC**$(x) = y$ *and* **OUT**$(x) \perp\!\!\!\perp $ **OUT**$(y) \vee($**IN**$(x)$ *and* **IN**$(y))$ **then**
        **LC**$(x) := NULL$
        Let $z$ be the smallest subset of **IN**$(x) \cup $ **IN**$(y)$ such that
        **OUT**$(x) \perp\!\!\!\perp $ **OUT**$(y) \vee (z)$
        **IN**$(x) := $ **IN**$(x) \cup z$
        **IN**$(y) := $ **IN**$(y) \cup z$
    **end**
**end**

**Algorithm 5:** Step 4 of the DM algorithm

29

Step 5. $pc.pattern.infinite := PC(\mathbf{X}, depth = infinite)$
Step 6. Examine the graphs produced in steps 4 and 5 (name these $G4$
and $G5$, respectively).
**for** *each output variable $O_i$ in $G4$ such that there is no direct edge*
*between $O_i$ and any input variables in $G5$* **do**
  remove the edge between $O_i$ and its latent.
  Add any adjacencies (from G5) between $O_i$ and the outputs
  connected to $O_i$'s former latent
**end**
Step 7. Return the graph from the end of step 6.
      **Algorithm 6:** Steps 5, 6, and 7 of the DM algorithm

## 3.4 Worked Example

Having described the algorithm (pages 28, 29, 29, 29, and 30), we can now move

on to a worked step-by-step example of the algorithm being used. Say we have

the following as the true graph (Figure 3.1):

Step 1: First, we run the PC algorithm with depth set equal to 0, giving

Figure 3.2.

Now we examine each variable, and if it has an an indegree of 0, we call it

an input. Otherwise, it is an output. In Figure 3.2, 1, 2, 3, and 4 are inputs,

while 5, 6, 7, 8, and 9 are all outputs.

Note; If the input variables are known, as is often the case, Step 1 can be

skipped (though input/output dependences still need to be found).

Step 2: Next, we look to see which output variables have a common set of

directly connected inputs. In this case, we have two sets: $OUT(< 5, 6, 7 >)$

is connected to $IN(< 1, 2 >)$, while $OUT(< 8, 9 >)$ is connected to $IN(<
1, 2, 3, 4 >)$. For each of these input/output pairs, we posit a latent variable

between the input and output .

Note: On the assumption that the system is linear, the number of latent

variables can be inferred from the rank of the correlation matrix of the out-

put variables – provided there are no causal connections among the measured
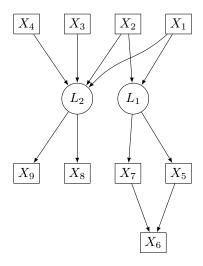
Figure 3.1: True graph

variables.

Step 3: We now check to determine which (if any) of the input sets are a subset of another input set. In this example, $IN(<1,2>)$ is a subset of $IN(<1,2,3,4>)$. We draw a path from the latent beneath $IN(<1,2>)$ to the latent beneath $IN(<1,2,3,4>)$. We also remove 1 and 2 from $IN(<1,2,3,4>)$, giving $IN(<3,4>)$. Finally, we draw a path from each input to the latent beneath the input, as well as draw a path to each output from the latent above the output.

This gives us the following graph (Figure 3.4):

Step 4: We now apply the next step in the algorithm (i.e., Sober (1998)). In this step, we check to see whether conditioning on some set of inputs (belonging to two latents connected by a directed edge) produces independence between the outputs of the two connected latents. If so, then we remove the path connecting the two latents, and draw paths from the conditioning set of inputs to each latent. In the example, this means that we remove the path between latents 1 and 2, and we connect $X_1$ and $X_2$ to both latents. Doing so yields the following
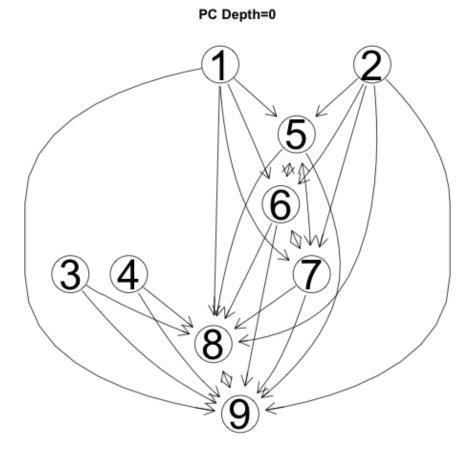
Figure 3.2: The pattern returned by the PC algorithm, with depth set to 0.



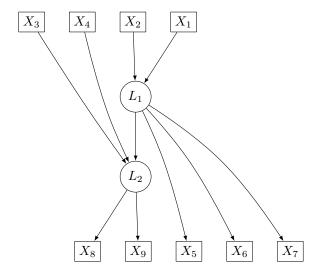Figure 3.3: List of all known variables.

Figure 3.4: The graph prior to applying Sober's criterion.

graph (Figure 3.5):

Step 5: Finally, we rerun the PC algorithm (but with a depth greater than zero). We continue running the PC algorithm (increasing the depth by one each time) until none of the reported graph adjacencies change.

Step 6: Using this graph (Figure 3.6), we look to see if any of the output variables have no direct edges connecting them to an output variable. If so, we then disconnect that (former) output variable from its latent, and connect the variable to the outputs (still connected to the latent) using the adjacencies reported in the $PC(depth = 0)$ graph. In the example, output 6 has ceased to have any direct edges connecting it to an input variable. Therefore, we disconnect $X_6$ from latent 1, and draw edges to $X_6$ from $X_5$ and $X_7$.

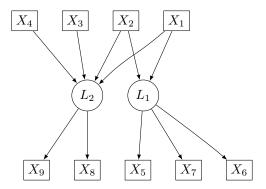Step 7: We now end the algorithm, and return the graph depicted in Figure 3.7.
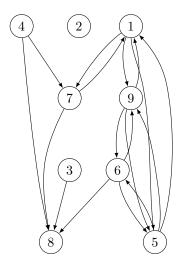
Figure 3.5: The graph after applying Sober's criterion.



Figure 3.6: The pattern reported by the PC algorithm, with depth greater than 0.
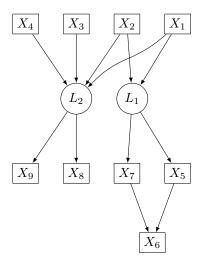
Figure 3.7: Final graph.

## 3.5    Proof of Correctness

**Proof of correctness for step 1 (classifying input and output variables)**

Due to assumptions A1, A2, and A3, the PC algorithm will produce a pattern consistent with the unconditional independence relations true of the measured variables in the true graph. Using this pattern, every input variable from the generating graph will only have adjacencies connecting it to output variables in the generating graph (as assumption A7 forbids adjacencies between input variables).

For every pair of variables that are inputs in the true graph, there will be no adjacency between the two variables in the PC pattern (by A7).

For every variable that is an input in the true graph and every output that is a descendant of that input, there will be an adjacency in the pattern returned by PC (by A6).

All of these adjacencies in the pc.pattern will ultimately be a directed edge from an input to an output variable, as the only paths from inputs to outputs in

35

the PC graph output will be through unshielded colliders (due to assumptions A5 and A6). Therefore, every input will have a total degree of no more than 0. Finally, due to assumption A8, every output variable must have an indegree greater than 0. So step 1 correctly classifies the input and output variables.

**Proof of correctness for step 2**

As every edge connecting an input to an output in pc.pattern must be the result of a path through a latent in the true graph (due to A6), and every output variable is a descendant of a latent (due to A8), there must be at least one latent (assuming the PC graph is not empty).

If there are sets of outputs whose members only have edges (in the pc.pattern) to some subset of the inputs, then there must be more than one latent (due to A6), and each of these sets of outputs must have its own latent as the only path from an input to an output is through a latent (again due to A6). This yields the correct number of latents.

**Proof of correctness for step 3**

If the input set of a latent (a) is a subset of the input set of another latent (b), and a is the largest such subset, then it must be the case that a is a latent cause of b (or latents a and b share some inputs). Otherwise, the inputs of a would have to have a path to the outputs of b via a non-latent (forbidden by A6), or via some latent between a and b (which is forbidden by the "largest subset" condition).

**Proof of correctness for step 4**

If step 3 reports an edge between two latents, then either that edge exists in the true graph, or the latents share some input variables (A4 forbids both being true simultaneously). Therefore, if there isn't an edge connecting the two latents in

the true graph, then $OUT(x) \perp\!\!\!\perp OUT(y) \vee (IN(x)IN(y))$, as there would be no open path connecting $OUT(L1)$ and $OUT(L2)$. If there is an edge between L1 and L2 in the true graph, then $OUT(x) \not\perp\!\!\!\perp OUT(y) \vee (IN(x)IN(y))$.

**Proof of correctness for step 5**

PC can be used due to A1, A2, and A3.

**Proof of correctness for step 6**

If an output variable has no paths to an input variable (in the pc.infinite pattern), then that output variable must be a child of only other output variables, else conditioning on observed variables would be insufficient to block all paths between the output variable and the input variables.

# Chapter 4

# Simulations

## 4.1  Why Perform a Simulation?

Despite the rotation problem (described in the literature review), factor analysis
remains a popular method for discovering causal structures. Analytic arguments
against a method are not always sufficiently persuasive on their own. Therefore,
simulations were run to support the analytic argument by providing both an
illustration of the weakness of factor analysis and a source of comparison for the
DM algorithm.

When searching for a MIMIC model structure, there are four main goals:
finding the correct number of latents, classifying variables into inputs and out-
puts, clustering inputs and outputs around their respective latent(s), and deter-
mining the latent-to-latent structure. As factor analysis is incapable of perform-
ing the second and fourth goals, and measures of accuracy are only meaningful
when being used as a source of comparison, the simulations only examine the
first and second.

## 4.2   How the Simulations were Performed

Multivariate Gaussian data were generated from seven different causal graphs
(see Figures 4.1 and 4.2 for a list of graphs). Each input variable followed a
standard normal distribution. As the latent variables were generated by a linear
combination of their inputs (and connected latents), they are simply a weighted
sum (where the edge weights are all set to 1) of standard normal variables. Noise
variables also followed a standard normal distribution. The outputs followed the
same distribution, as they too were a weighted sum of standard normal variables.
Factor analysis and the DM algorithm were both run on the same datasets, with
sample sizes of 250, 500, 1000, and 10,000 observations. For each causal graph
and sample size, the data were generated five hundred times and the algorithms
run and scored on each dataset. Setting the edge weights to 1 is convenient but
not optimal to demonstrate the accuracies of various algorithms, but random
assignment of weights would only change the strengths of covariances, not the
asymptotic constraints on the data that various models (e.g., a 1 factor model)
imply.

As it is recommended to select the number of latents in factor analysis us-
ing several different methods (Hair et al., 1998), four different methods (non-
graphical approximations of a scree plot) were used: optimal coordinates, accel-
eration factor, parallel analysis, and the Kaiser rule. Each method was allowed
a vote for the recommended number of factors, and the number with the highest
frequency was selected. In the event of a tie, the smallest number was selected.

The algorithms were judged based on two criteria: the number of cases where
a given algorithm reported an incorrect number of latents (reported as the per-
centage of incorrect cases out of 500), and (assuming the algorithm reported the
correct number of latents) the reliability of undirected edge discovery. Relia-

39

bility was represented by the true positive[1], false positive[2], and true discovery[3] rates for reported edges. Reported graphs were judged on their ability to identify undirected edges correctly, as it is unclear how without prior knowledge measured variables are to be separated into input (exogenous) variables and output (endogenous) in a reported factor model.

---

[1]The true positive rate is the number of correctly found edges divided by number of true edges in the actual graph.

[2]The false positive rate is the number of incorrectly found edges divided by the number of true gaps in the actual graph.

[3]The true discovery rate is the number of correctly found edges divided by number of found edges.
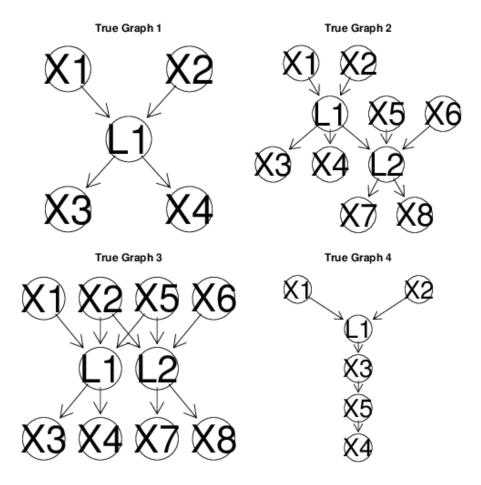
Figure 4.1: The first four graphs tested. Note that L1 and L2 are latent variables.
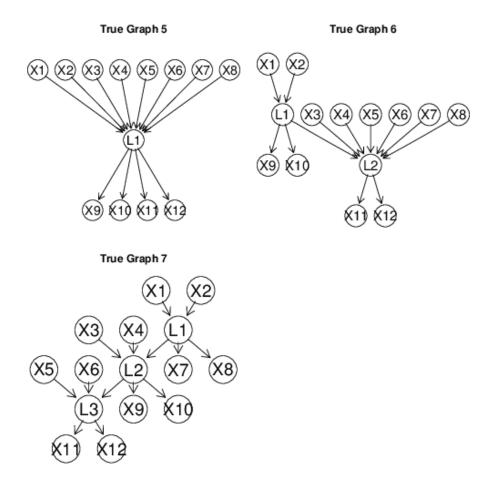
Figure 4.2: The second three graphs tested. Note that L1, L2, and L3 are latent variables.

## 4.3  Analysis of Results

In general, the DM algorithm outperformed factor analysis, with a lower false positive rate (graphs 2, 3, 4, and 6), a higher true discovery rate (graphs 2, 3, 4, 6, and 7), and a higher true positive rate (graphs 2, 3, 4, 6, and 7). The algorithm also proved to be correct about the number of latents more often, with the exception of graphs 3 and 5 where it performed marginally worse, and graph 6, where there were relatively many variables and few latents (as a single incorrect independence test can result in the creation of a new latent, this increased the odds of getting the number of latents wrong). Performance on graph 1 was similar for both methods. Finally, unlike factor analysis, the performance of the DM algorithm generally improved as sample size increased. Bar plots of accuracies are given in Appendix A.

Factor analysis is unreliable at best. In some cases, its performance decreases as sample size increases! For graphs 2 and 7, the number of correct latent cases became smaller as sample size increased (in some cases, factor analytical methods did not get the right answer once). In most cases, factor analysis was outperformed by the DM algorithm, except for graphs 1 and 5, where it marginally performed better. However, given the general unreliability of factor analysis for correctly determining even a small number of latents, such a performance cannot be reasonably extended to general cases (such as those encountered during a data analysis).

# Chapter 5

# Data Analysis

## 5.1   Goals of the Analysis

The two primary goals of analysis were:

1. Cluster the genes (inputs) and expressions (outputs) around their respective latents.

2. Find latent-to-latent connections.

As a secondary goal, any unusual patterns, such as an uncommon or atypical number of input variables tied to a single latent in the resulting MIMIC graph, are noted for future study.

## 5.2   Algorithm Assumptions and the Dataset

In the dataset, many of the algorithm's assumptions are known to be satisfied. Inputs are independent of one another (as genes do not "cause" one another), no input has a path to an output except via a latent (it is impossible for a gene

to directly affect mRNA), and every measured variable is either an input or an output.

As variables are already classified as inputs or outputs, we can relax the assumption that every latent have 2 inputs and 2 outputs. Now they need only have 1 of each. This also means that we do not need to run the PC algorithm, but can instead simply record which inputs and outputs are dependent. Since it is impossible for mRNA (output) variables to directly cause other mRNA variables, we can also omit the other run of PC used in the last step of the DM algorithm.

The true graph may prove to not be singly connected, however this only means that some information will be lost (e.g., an input shared by 2 connected latents may be reported as being directly connected to only 1 latent). In which case, the algorithm is reporting one of several graphs consistent with the observed independence relations.

The most dubious assumption made is that of an acyclic graph. Cyclic relations between gene expressions and latent variables have been observed in other genomic datasets. In the event of this assumption being violated, latents can be incorrectly merged, and in some cases the reported latent-to-latent structure can be incorrect. At present there is no known solution to this problem, and any results reported by the algorithm should be interpreted with this caveat in mind.

## 5.3   The Initial Data

The original data were gathered using microarray analyses and massively parallel sequencing (Network et al., 2011). Every observation was a patient with ovarian cancer. Each patient had two classes of variables recorded. First, every

gene was categorized based on whether it was mutated[1] or had multiple copies (i.e., more than a single possibly partial copy of a chromosome) in some fashion. If so, that gene was coded as "1". Otherwise, it was coded as "0".

Second, gene expression levels, in the form of mRNA, were measured as continuous variables. This continuous data was then split into three-level ordinal variables. The cutoffs[2] were chosen based on 25% and 75% of the range of values, with below 25% coded as "-1" (for low gene expression level), above "75%" as "1" (for high), and all others cases as "0" (for normal).

At the beginning of analysis, the dataset consisted of 562 patients, all with ovarian cancer. There were 17,610 genes (inputs), as well as 12,042 gene expression variables (outputs).

The mean number of mutations was 205 and the median was 124.5. The largest number of mutations any subject had was 2,968. There were no subjects with zero abnormalities. Figure 5.1 depicts the frequency distribution of abnormalities per subject.

The average gene was mutated in 6.544 subjects, though the median number of subjects was 4. The most frequently mutated gene was mutated in 382 subjects. Figure 5.2 depicts their distribution.

---

[1]Meaning a change in at least a base or an allele due to: non-synonymous substitution, frameshift, deletion, or truncation.
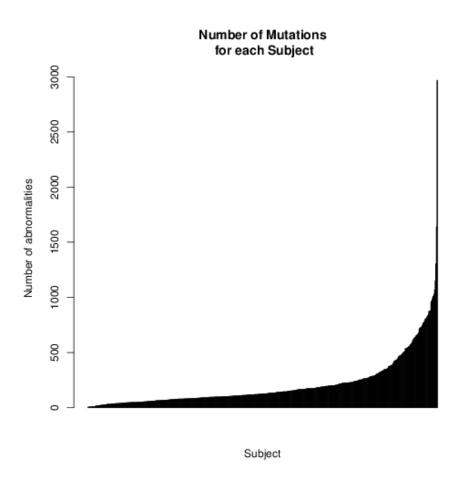
[2]The cutoffs were chosen by the creator of the dataset.

**Number of Mutations
for each Subject**

Figure 5.1: Number of abnormalities for each subject.

47

**Mutation Frequency for each Gene**

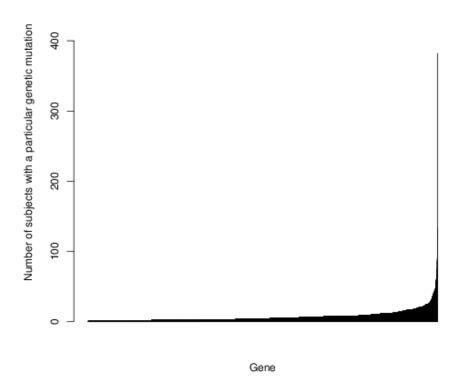*Number of subjects with a particular genetic mutation*

*Gene*

Figure 5.2: Number of subjects with a given abnormality.

## 5.4 Dimension Reduction

Due to computational constraints in both memory and time, the algorithm cannot (as implemented) be run on the complete dataset. Therefore, cross-validated[3] lasso regression[4] (Hastie et al., 2009) was used in order to reduce the number of variables in the dataset. This regression was used in two ways. First, it was used to remove gene variables by predicting expression levels using genes, taking advantage of the sparsity forced by the lasso regression to exclude less useful gene variables. Similarly, lasso regression was used to remove expression variables by predicting genes using expression levels, again exploiting the forced sparsity property of lasso regression. Lasso regression was used (instead of simply randomly selecting some number of variables for inclusion) in order to increase the chance a selected node was related to another node.

As a constant cannot be used as the response of a regression, constant gene and expression variables were dropped.

The reason categorical expression variables were used instead of continuous versions is that the function used to perform the cross-validated lasso regression reported ambiguous errors (likely due to some variables having very few instances of one of the possible categories) when using a binomial response (i.e., when predicting whether or not a gene was mutated). This problem of ambiguous errors necessitated treating the gene variables as continuous. To avoid mixing variable types, the expression variables were also treated as continuous, despite being ordinal categoricals. Doing this had the added benefit of increasing the rate with which the lasso regressions could be performed.

Even after using lasso regression to select variables, too many variables re-

---

[3]Mean squared error was used as the criterion for judging models.

[4]As a number of variables had some categories which rarely occurred, the lasso regression assumed the variables were continuous. Doing otherwise would have resulted in many variables being dropped for no reason other than a lack of variability. Note that as a result, no regression coefficients produced by the lasso are interpretable.

mained. Therefore, an additional reduction was performed. For each of the 29,652 lasso regressions (17,610 with genes as the response, 12,042 with expressions as the response), the number of chosen predictors was examined. Only variables that belong to an unusually large group of predictors (of a size greater than the 99th quantile) were preserved. In other words, a variable was only included if it was part of a large collection of predictors in at least one lasso regression. Specifically, a group had to be greater than 69 predictors for gene variables and 141 predictors for expression variables in order for its members to be preserved in the final dataset. This reduced the total number of variables in the dataset to 4,369.

### 5.4.1 Side Effects of Dimension Reduction

As any dimension reduction procedure necessarily results in the exclusion of some variables, violations of the various assumptions made by the DM algorithm can occur no matter what reduction procedure is followed. The algorithm (when used on the cancer dataset) assumes that each latent variable has at least one gene mutation (input) and gene expression (output) variable. If the dimension reduction procedure were to drop all of a latent's inputs or outputs, then the structure reported by the algorithm can be incorrect. As a latent might only posses a single input or output, simply excluding 1 variable can lead to errors. One possible error, which may partially explain the lack of latent-to-latent edges in the graphs shown in Figures 5.4 and 5.5 (though such absences could also be the truth), occurs when a latent losses all of its outputs. Consider the following graph (Figure 5.3). If the algorithm were run on the observed variables in it, the resulting graph would have no edges between latents, even though the true graph contains such an edge (It would also report too few latents). As we do not already know the underlying causal structure in the dataset, we
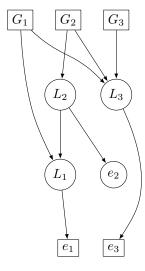
Figure 5.3: The above graph gives an example where, due to a dimension reduction procedure, the output variable $e_2$ has been omitted, resulting in the DM algorithm returning an incorrect graph.

cannot determine if this situation has occurred. It is thus advisable to avoid any dimension reduction when deploying the DM algorithm. Due to the limited computational resources available, this guideline had to be ignored. Caveat lector.

## 5.5    Graphical Results

To construct the initial adjacency matrix, Fisher's exact test (of independence) was used instead of the PC algorithm, as inputs and outputs were already known and doing so reduced the number of independence tests performed. For Sober's step, a chi-squared test was used. A p-value of .000001 was used when performing both kinds of tests.

After running the algorithm on the reduced dataset and removing all nodes with a total degree of zero, 445 variables remained. The reported graph (Fig-

**Protein Signalling Network**



Figure 5.4: Initial reported signalling network. Note that due to the number of nodes, edges can sometimes overlap.

ure 5.4) contains 139 latents, 295 inputs, and 150 outputs.

### 5.5.1 Interpretation

Interestingly, there are a number of disconnected subgraphs in the reported graph (Figure 5.4). There are three subgraphs where a single output variable is related to a large cluster of input variables. Although this is a possible consequence of the dimension reduction procedure (as relevant gene expression level

variables capable of "breaking up" the cluster may have been excluded), the result may point to gene expression variables usable as indicators for the presence of ovarian cancer. There are also several subgraphs with more complicated network structures. Enlarged graphs of these clusters are included in Appendix B.

One interesting thing to note (in both Figures 5.4 and 5.5) is the relative absence of latent-to-latent connections (indeed, there was only a single latent-to-latent edge. Latent 88 caused latent 137 in Figure B.22). One possible explanation (covered in more detail in the section on dimension reduction) for this phenomenon is that, due to the dimension reduction procedure, a single latent is "standing in" for a number of others. Therefore, while there may be many latent-to-latent connections, their existence is being covered up by the presence of some disconnected latents that are wrapped up in the single "representative" latent. An alternative explanation is that the p-value chosen for performing Sober's step was too extreme.

As a nice sanity check for the reported graph, note that several gene variables (HOOK2, GIPC1, NDUFB7, and PIN1) have directed paths to their (known to be) related gene expressions (HOOK2.1, GIPC1.1, NDUFB7.1, and PIN1.1). An unsystematic search of the literature failed to discover any clear association between regulatory role and latent clustering.

**Protein Signalling Network**

Figure 5.5: A subgraph of the larger protein signalling network. Note that the black numbers are simply labels for the latent variables. Also note that some output variables have a ".1" appended to their name in order to avoid naming conflicts with input variables. There is an enlarged version in Appendix B.

# Chapter 6

# Future Work

There are a number of different possible avenues for future research.

The implementation using backwards and forwards regression for variable elimination and a very low alpha value probably resulted in a structure for the genomics data that is too sparse in important respects. In a practical scientific application, the algorithm would need to be run with variations in these strategies, and/or using a false discovery rate. The algorithm could be improved if methods can be found that allow relaxing various assumptions (BPC and MIMBuild seem likely candidates), most importantly:

1. Acyclicality (A3).

2. Inputs are independent of one another (A7).

3. Singly connected graph (A4).

Some properties of the algorithm, such as its computational complexity could also be formally calculated. Another possible avenue for research is more thorough simulations, involving:

1. Randomly generated MIMIC models, with different numbers of latents, input-to-output variable ratios, and latent-to-observable variable ratios.

2. Different parameterizations (and distributions).

It is also worthwhile to determine how robust the results of the algorithm are given different dimension reduction (and other data analysis) decisions. For example, how does the reported graph change when the p-value used to perform Sober's step changes?

In the event more computational power becomes available, running the algorithm on the entire genomic dataset would be of interest, as well as testing the implications of the results. As the computational power needed to run the algorithm on a dataset of around 30,000 variables is very high, a parallelized version of the algorithm, as well as access to a supercomputer may need to be investigated. There are also a number of other datasets, including those from neuroscience, psychology, economics, and other subfields of genomics, on which the algorithm could be run.

# Appendix A

# Simulation Results

This appendix contains the results for the simulations discussed in chapter 4.

Figure A.1: The first four graphs tested. Note that L1 and L2 are latent variables.

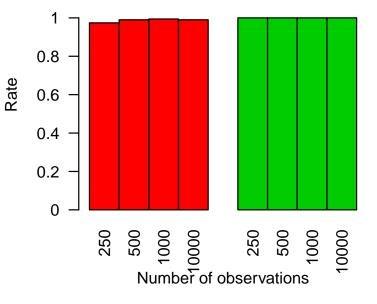Figure A.2: The second three graphs tested. Note that L1, L2, and L3 are latent variables.
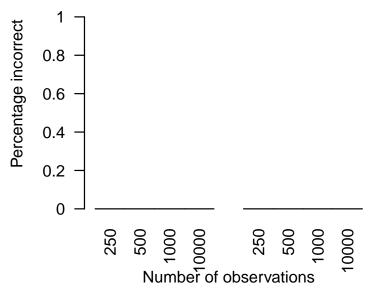
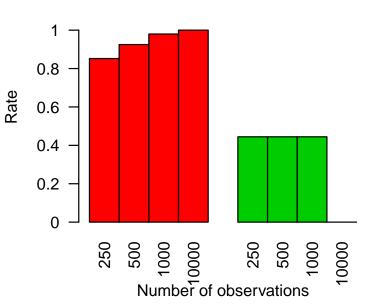# Graph 1 True Positive Rate

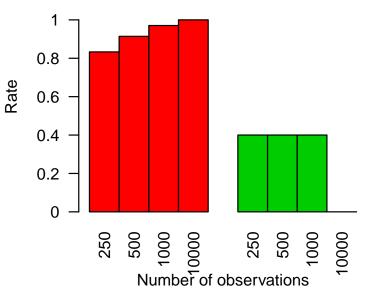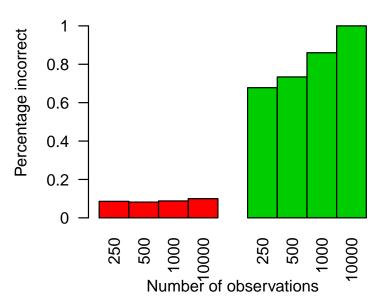# Graph 1 False Positive Rate

# Graph 1 True Discovery Rate

# Graph 1 Precentage of False Latent Cases (out of 500)

60
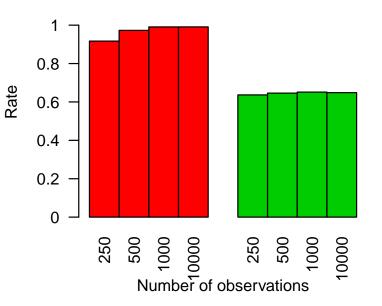
# Graph 4 True Positive Rate

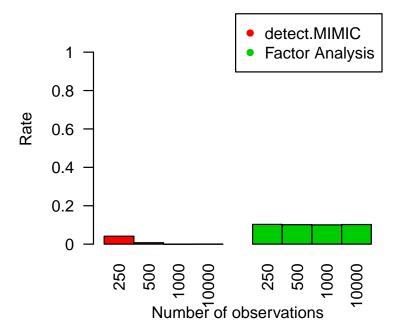# Graph 4 False Positive Rate

# Graph 4 True Discovery Rate
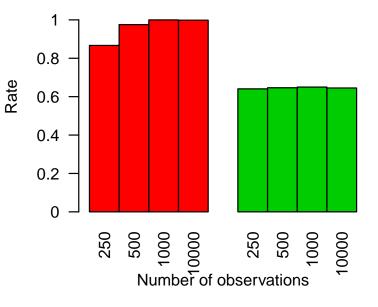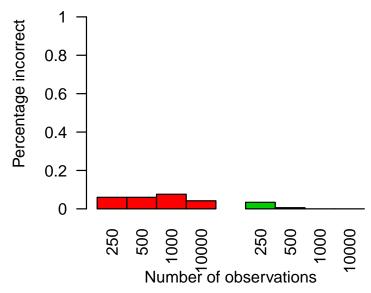
# Graph 4 Precentage of False
# Latent Cases (out of 500)

**Graph 6 True Positive Rate**

**Graph 6 False Positive Rate**

**Graph 6 True Discovery Rate**

**Graph 6 Precentage of False Latent Cases (out of 500)**

# Appendix B

# Data Analysis

This appendix contains both an enlarged version of the full protein signalling network, as well as enlarged versions of every subgraph with more than 3 nodes.

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network



71

# Protein Signalling Network



72

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Protein Signalling Network

# Appendix C

# Souce Code

This appendix contains the source code used when performing both the simulations and the data analysis.

```r
1   require(pcalg)
2   require(igraph)
3   require(gRim)
4   require(plotrix)
5   require(nFactors)
6
7
8
9   source("construct_graph.R")
10  source("data_commands.R")
11  source("simple_simulations.R")
12  source("simulation_scoring.R")
13
14
15  # Currently, this portion is not in a function so as to avoid rerunning
16  # everything in the event of an error (i.e., an error will not interrupt the
17  # function call, leading to already completed simulations being discarded).
18  # Eventually, a function will be written so that this is cleaner.
19
20  # Runs the simulations.
21  graph1.10000 <- replicate(test.both.methods("sim.graph.1.r.txt",
22    sample.size=10000), n=500)
23  graph2.10000 <- replicate(test.both.methods("sim.graph.2.r.txt",
24    sample.size=10000), n=500)
25  graph3.10000 <- replicate(test.both.methods("sim.graph.3.r.txt",
26    sample.size=10000), n=500)
27  graph4.10000 <- replicate(test.both.methods("sim.graph.4.r.txt",
28    sample.size=10000), n=500)
29  graph5.10000 <- replicate(test.both.methods("sim.graph.5.r.txt",
30    sample.size=10000), n=500)
31  graph6.10000 <- replicate(test.both.methods("sim.graph.6.r.txt",
32    sample.size=10000), n=500)
33  graph7.10000 <- replicate(test.both.methods("sim.graph.7.r.txt",
34    sample.size=10000), n=500)
35
36
37  graph1.1000 <- replicate(test.both.methods("sim.graph.1.r.txt",
38    sample.size=1000), n=500)
39  graph2.1000 <- replicate(test.both.methods("sim.graph.2.r.txt",
40    sample.size=1000), n=500)
41  graph3.1000 <- replicate(test.both.methods("sim.graph.3.r.txt",
42    sample.size=1000), n=500)
43  graph4.1000 <- replicate(test.both.methods("sim.graph.4.r.txt",
44    sample.size=1000), n=500)
45  graph5.1000 <- replicate(test.both.methods("sim.graph.5.r.txt",
46    sample.size=1000), n=500)
47  graph6.1000 <- replicate(test.both.methods("sim.graph.6.r.txt",
48    sample.size=1000), n=500)
49  graph7.1000 <- replicate(test.both.methods("sim.graph.7.r.txt",
50    sample.size=1000), n=500)
51
52  graph1.500 <- replicate(test.both.methods("sim.graph.1.r.txt",
53    sample.size=500), n=500)
54  graph2.500 <- replicate(test.both.methods("sim.graph.2.r.txt",
55    sample.size=500), n=500)
56  graph3.500 <- replicate(test.both.methods("sim.graph.3.r.txt",
57    sample.size=500), n=500)
58  graph4.500 <- replicate(test.both.methods("sim.graph.4.r.txt",
59    sample.size=500), n=500)
60  graph5.500 <- replicate(test.both.methods("sim.graph.5.r.txt",
61    sample.size=500), n=500)
62  graph6.500 <- replicate(test.both.methods("sim.graph.6.r.txt",
63    sample.size=500), n=500)
64  graph7.500 <- replicate(test.both.methods("sim.graph.7.r.txt",
65    sample.size=500), n=500)
66
67  graph1.250 <- replicate(test.both.methods("sim.graph.1.r.txt",
```

91

```
68      sample.size=250), n=500)
69   graph2.250 <- replicate(test.both.methods("sim.graph.2.r.txt",
70      sample.size=250), n=500)
71   graph3.250 <- replicate(test.both.methods("sim.graph.3.r.txt",
72      sample.size=250), n=500)
73   graph4.250 <- replicate(test.both.methods("sim.graph.4.r.txt",
74      sample.size=250), n=500)
75   graph5.250 <- replicate(test.both.methods("sim.graph.5.r.txt",
76      sample.size=250), n=500)
77   graph6.250 <- replicate(test.both.methods("sim.graph.6.r.txt",
78      sample.size=250), n=500)
79   graph7.250 <- replicate(test.both.methods("sim.graph.7.r.txt",
80      sample.size=250), n=500)
81
82
83   save.image("simulation.results.latents.RData")
84
85
86
87
88   score.graph1.10000<- score.both(graph1.10000, graph.name="sim.graph.1.r.txt")
89   score.graph1.1000<- score.both(graph1.1000, graph.name="sim.graph.1.r.txt")
90   score.graph1.500<- score.both(graph1.500, graph.name="sim.graph.1.r.txt")
91   score.graph1.250<- score.both(graph1.250, graph.name="sim.graph.1.r.txt")
92
93   score.graph2.10000<- score.both(graph2.10000, graph.name="sim.graph.2.r.txt")
94   score.graph2.1000<- score.both(graph2.1000, graph.name="sim.graph.2.r.txt")
95   score.graph2.500<- score.both(graph2.500, graph.name="sim.graph.2.r.txt")
96   score.graph2.250<- score.both(graph2.250, graph.name="sim.graph.2.r.txt")
97
98   score.graph3.10000<- score.both(graph3.10000, graph.name="sim.graph.3.r.txt")
99   score.graph3.1000<- score.both(graph3.1000, graph.name="sim.graph.3.r.txt")
100  score.graph3.500<- score.both(graph3.500, graph.name="sim.graph.3.r.txt")
101  score.graph3.250<- score.both(graph3.250, graph.name="sim.graph.3.r.txt")
102
103  score.graph4.10000<- score.both(graph4.10000, graph.name="sim.graph.4.r.txt")
104  score.graph4.1000<- score.both(graph4.1000, graph.name="sim.graph.4.r.txt")
105  score.graph4.500<- score.both(graph4.500, graph.name="sim.graph.4.r.txt")
106  score.graph4.250<- score.both(graph4.250, graph.name="sim.graph.4.r.txt")
107
108  score.graph5.10000<- score.both(graph5.10000, graph.name="sim.graph.5.r.txt")
109  score.graph5.1000<- score.both(graph5.1000, graph.name="sim.graph.5.r.txt")
110  score.graph5.500<- score.both(graph5.500, graph.name="sim.graph.5.r.txt")
111  score.graph5.250<- score.both(graph5.250, graph.name="sim.graph.5.r.txt")
112
113  score.graph6.10000<- score.both(graph6.10000, graph.name="sim.graph.6.r.txt")
114  score.graph6.1000<- score.both(graph6.1000, graph.name="sim.graph.6.r.txt")
115  score.graph6.500<- score.both(graph6.500, graph.name="sim.graph.6.r.txt")
116  score.graph6.250<- score.both(graph6.250, graph.name="sim.graph.6.r.txt")
117
118  score.graph7.10000<- score.both(graph7.10000, graph.name="sim.graph.7.r.txt")
119  score.graph7.1000<- score.both(graph7.1000, graph.name="sim.graph.7.r.txt")
120  score.graph7.500<- score.both(graph7.500, graph.name="sim.graph.7.r.txt")
121  score.graph7.250<- score.both(graph7.250, graph.name="sim.graph.7.r.txt")
122
123
124
125
126  # TODO: Rewrite this portion. It is ugly, and should be in a function.
127  number.graphs<-7
128  graph.groups <- c()
129  for(k in 1:number.graphs){
130      graph.groups[[k]] <- list(
131       ls()[grep(pattern=paste("score.graph",
132       k, ".", sep=""), ls(), fixed=T)])
133  }
134
```

```
135
136    save.image("scored.simulations.latents.RData")
137
138
139
140    pdf("sim.res_final.pdf")
141    plot.rates(convert.graph.groups(graph.groups))
142    dev.off()
```

```r
 1
 2    # converts the Tetrad .r.txt representation to a graphNEL object.
 3    read.dag <- function(file){
 4        orig.mat <- read.table(file=file)
 5
 6        orig.mat[orig.mat==1] <-0
 7        orig.mat[orig.mat==-1] <-1
 8
 9
10        final.graph <- igraph.to.graphNEL(graph.adjacency(as.matrix(orig.mat)))
11
12        return(final.graph)
13    }
14
15
16    # Takes in a graphNEL object, and generates normally distributed data from it.
17    generate.data.from.dag <- function(graph, n=100, errDist="normal"){
18
19        top.sort <- topological.sort(igraph.from.graphNEL(graph))
20        var.names <- nodes(graph)[top.sort]
21
22        graph<-igraph.to.graphNEL(graph.adjacency(as(graph, "matrix")[var.names,
23         var.names]))
24
25        generated.data<-data.frame(rmvDAG(dag=graph, n=n, errDist=errDist))
26        names(generated.data) <- var.names
27
28        return(generated.data)
29    }
30
```

```r
1
2    # Runs both FA test, as well as detect.MIMIC
3    test.both.methods <- function(file="sim.graph.1.r.txt", sample.size=1000,
4     alpha=.01, pval=.05, cut.off=.3, scree=FALSE){
5
6        dataset <- generate.data.set(file=file, sample.size=sample.size)[[1]]
7
8        clean.dataset <- dataset[,-grep(pattern="[L:digit:]",
9         names(dataset))]
10
11       # print(head(clean.dataset))
12
13       fa.results <- test.fa(clean.dataset,
14            n.latents=ncol(dataset)-ncol(clean.dataset), cut.off=cut.off,
15            scree==scree)
16
17       detect.mimic.results <- find.mimic(clean.dataset, alpha=alpha,
18               pval=pval)
19
20
21       return(list(fa.results=fa.results,
22            detect.mimic.results=detect.mimic.results))
23
24    }
25
26    # Runs runs factor analysis on datset. If scree=TRUE, then attempts to
27    # discover n.latents. Uses varimax rotation.
28    test.fa <- function(dataset, n.latents=1, cut.off=.3, scree=TRUE){
29
30        var.names <- names(dataset)
31
32        if(scree){
33            # As most papers have suggested using several different tests for
34            # n.factors, have chosen to choose most frequent number reported. In
35            # the event of a tie, the smaller number of latents is chosen, as
36            # most papers suggest parsimoney.
37            n.latents <-
38             as.numeric(names(which.max(table(unlist(c(nScree(x=dataset,
39             model="factors")$Components))))))
40            fa.model <- factanal(dataset, factors=n.latents)
41        }
42        else{
43            fa.model <- factanal(dataset, factors=n.latents)
44        }
45
46
47        return(fa.model)
48    }
49
50
51
52    generate.data.set <- function(file="sim.graph.1.r.txt", sample.size=10000){
53        graph.data <- list(generate.data.from.dag(read.dag(file),
54         n=sample.size))
55
56        return(graph.data)
57    }
58
59
60    test.graphs <- function(n.graphs=1:7, seed.set=NULL, sample.size=1000, plot.graphs=TRUE){
61
62        if(!is.null(seed.set)){set.seed(seed.set)}
63                                                    95
64        found.graphs <- list()
65
66        files.to.load<-paste("sim.graph.",n.graphs,".r.txt", sep="")
67        for(i in files.to.load){
```

```r
68
69              current.graphs <- test.graph(i, sample.size=sample.size)
70
71              found.graphs<-list(found.graphs, current.graphs)
72
73              if(plot.graphs){plot.test(current.graphs)}
74          }
75      return(found.graphs)
76  }
77
78  # Reads in graph file, and generates normally distributed data from it.
79  test.graph <- function(file, sample.size=1000){
80      true.graph<-read.dag(file=file)
81
82      generated.data <- generate.data.from.dag(graph=true.graph, n=sample.size)
83
84      generated.data.no.latents <- generated.data[,-grep(pattern="[L:digit:]",
85       names(generated.data))]
86
87      result<-find.mimic(generated.data.no.latents)
88      return(list(result=result, true.graph=true.graph))
89  }
90
91  # Plots the various stages of the algorithm used in finding the graph.
92  plot.test<-function(graph.list){
93
94      results<-graph.list$result
95      true.graph<-graph.list$true.graph
96
97      # No results, so just plots true graph
98      if(is.null(results)){plot(true.graph); return()}
99      else if(class(results)!="list"){plot(results)}
100     else{
101         print(results)
102
103         par(mfrow=c(2,3))
104         if(!is.null(true.graph)){plot(true.graph, main="True Graph")}
105
106         if(!is.null(results$pc.depth.0)){plot(results$pc.depth.0,
107             main="PC Depth=0")}
108
109         if(!is.null(results$pre.sober.model)){
110             plot(results$pre.sober.model, main="Pre-Sober")}
111
112         if(!is.null(results$mimic.model.graph)){
113             plot(results$mimic.model.graph, main="Post-Sober")}
114
115         if(!is.null(results$last.pc)){plot(results$last.pc,
116             main="PC Depth>0")}
117
118         if(!is.null(results$final.model)){plot(results$final.model,
119             main="Final Graph")}
120     }
121
122  }
123
```

```r
1    # Returns average fa.score and mimic.score. Note that scores are
2    # as follows:
3
4    #           tpr
5    #           True Positive Rate: Number of correctly found edges (in estimated
6    # graph) divided by number of true edges (in true graph)
7    #
8    #           fpr
9    #           False Positive Rate: Number of incorrectly found edges divided by
10   #   number of true gaps (in true graph)
11   #
12   #           tdr
13   #           True Discovery Rate: Number of correctly found edges divided by
14   #   number of found edges (both in estimated graph)
15   score.both <- function(graph.sim, graph.name, cut.off=.3){
16
17       true.graph <- as(read.dag(graph.name), "matrix")
18
19       var.names <- names(data.frame(true.graph))
20
21       latent.positions <- grep(pattern="[L:digit:]", var.names)
22
23       non.latent.positions <- (1:length(var.names))[-latent.positions]
24
25       sorted.truth <- (1:length(var.names))[c(non.latent.positions,
26           latent.positions)]
27
28       true.graph <- true.graph[sorted.truth,sorted.truth]
29
30       # Scores FA model.
31       fa.score <- lapply(graph.sim[1,],
32           function(fa.model){score.fa(fa.model=fa.model, cut.off=cut.off,
33               true.graph=true.graph)})
34
35       # Scores MIMC model.
36       mimic.score <- lapply(graph.sim[2,],
37           function(mimic.model){
38               # Handles case when only PC output was returned.
39               if(length(mimic.model)>1){
40                   result<-score.mimic(mimic.model=mimic.model$final.model,
41                true.graph=true.graph)
42
43                   return(result)
44                   }
45               return(c(0,0,0))
46               })
47
48       n.with.incorrect.latents.mimic <- 0
49       for(i in mimic.score){
50           if(is.null(i)){n.with.incorrect.latents.mimic <-
51               n.with.incorrect.latents.mimic+1}
52       }
53
54       n.with.incorrect.latents.fa <- 0
55       for(i in fa.score){
56           if(is.null(i)){n.with.incorrect.latents.fa <-
57               n.with.incorrect.latents.fa+1}
58       }
59
60       fa.score <- (do.call(rbind, fa.score))
61       mimic.score <- (do.call(rbind, mimic.score))
62
63       if(is.null(fa.score)){fa.score <- c(0,0,0)}97
64       else{fa.score <- colMeans(fa.score)}
65       mimic.score <- colMeans(mimic.score)
66
67       return(list(fa.score=fa.score, mimic.score=mimic.score,
```

```R
68              latent.incorrect.mimic = n.with.incorrect.latents.mimic,
69              latent.incorrect.fa = n.with.incorrect.latents.fa))
70  }
71
72  score.mimic <- function(mimic.model, true.graph){
73
74      adj.matrix.true <- data.frame(true.graph)
75
76      adj.matrix.mimic <- data.frame(as(mimic.model, "matrix"))
77
78      n.latents.truth <- grep(pattern="[L:digit:]",
79       names(adj.matrix.true))
80
81      n.latents.mimic <- grep(pattern="[L:digit:]",
82       names(adj.matrix.mimic))
83
84      true.graph<-igraph.to.graphNEL(graph.adjacency(true.graph))
85
86
87      # If the mimic.model found has a different numnber of latents than the
88      # true graph, then cannot calculate tpr, fpr, tdr. Have therfore treated
89      # those as NULL objects. (i.e., as with FA, correct n.latents is assumed)
90      if(length(nodes(mimic.model)) == length(nodes(true.graph))){
91          return(compareGraphs(mimic.model, true.graph))
92      }
93      return(NULL)
94  }
95
96  # Scores FA model.
97  score.fa <- function(fa.model, cut.off=.3, true.graph){
98
99      fa.mat <- as(fa.model$loadings, "matrix")
100     n.latents <- ncol(fa.mat)
101     var.names <- row.names(fa.mat)
102     n.vars <- nrow(fa.mat)
103
104     fa.model <- prune.fa.paths(fa.model, cut.off=cut.off)
105
106     fa.model <- igraph.to.graphNEL(graph.adjacency(fa.model))
107     true.graph <- igraph.to.graphNEL(graph.adjacency(true.graph))
108
109     if(length(nodes(fa.model))==length(nodes(true.graph))){
110
111         graph.comparison <- (compareGraphs(fa.model, true.graph))
112         return(graph.comparison)
113     }
114     return(NULL)
115 }
116
117
118 prune.fa.paths <- function(fa.model, cut.off=.3){
119
120     fa.loadings.matrix <- as(fa.model$loadings, "matrix")
121
122     n.latents <- ncol(fa.loadings.matrix)
123     var.names <- row.names(fa.loadings.matrix)
124     n.vars <- length(var.names)
125
126     adj.mat <- matrix(FALSE, nrow=(n.vars+n.latents),
127      ncol=(n.vars+n.latents), dimnames=list("row"=c(var.names,
128         1:n.latents), "col"=c(var.names, 1:n.latents)))
129
130     for(i in 1:n.latents){
131
132         adj.mat[abs(fa.model$loadings[,i])>cut.off,
133         i+length(var.names)] <- TRUE
134     }
```

98

```R
135
136        adj.mat[(length(var.names)+1):(length(var.names)+n.latents),
137         (length(var.names)+1):(length(var.names)+n.latents)] <- FALSE
138
139        return(adj.mat)
140
141    }
142
143    get.latent.cluster <- function(adj.matrix, n.latents, n.vars){
144        latent.vectors.col <- adj.matrix[,(n.vars+1):(n.vars+n.latents)]
145        latent.vectors.row <- t(adj.matrix[(n.vars+1):(n.vars+n.latents),])
146
147        latent.clusters <- (latent.vectors.col+latent.vectors.row)[
148        -((n.vars+1):(n.vars+n.latents)),]
149        return(latent.clusters)
150    }
151
152
153    convert.graph.groups <- function(graph.groups = graph.groups){
154        graph.final<-list()
155
156
157        for(i in 1:length(graph.groups)){
158            graph.list <- unlist(graph.groups[[i]])[c(3,4,1,2)]
159            n.null.mimic <- c()
160            n.null.fa <- c()
161            graph.fa.score <- c()
162            graph.mimic.score <- c()
163            var.names <- c()
164            for(j in 1:length(graph.list)){
165
166                graph.score<- get(graph.list[j])
167
168                graph.fa.score <- rbind(graph.fa.score, unlist(graph.score[[1]]))
169                graph.mimic.score <- rbind(graph.mimic.score,
170                    unlist(graph.score[[2]]))
171
172                #   TODO: need to change this bit so that it handles the two-null
173                # case.
174                n.null.mimic <- c(n.null.mimic, unlist(graph.score[[3]]))
175                n.null.fa <- c(n.null.fa, unlist(graph.score[[4]]))
176            }
177
178            row.names(graph.fa.score) <- c("250", "500", "1000", "10000")
179            row.names(graph.mimic.score) <- c("250", "500", "1000", "10000")
180
181            graph.final[[i]] <- list(fa.scores=graph.fa.score,
182                mimic.scores=graph.mimic.score, n.null.mimic=n.null.mimic,
183                n.null.fa=n.null.fa)
184
185        }
186        return(graph.final)
187    }
188
189    plot.rates <-function(score.list){
190            for(i in 1:length(score.list)){
191
192                par(mfrow=c(2,2))
193
194                barplot(cbind(score.list[[i]]$mimic.scores[,1],
195                    score.list[[i]]$fa.scores[,1]),
196                   main=paste("Graph ",i, " True Positive Rate", sep=""),
197                    col=c(2,2,2,2, 3,3,3,3),       99
198                    ylab="Rate", xlab="Number of observations",
199                    ylim=c(0,1.2), beside=T, names=c("250", "500", "1000",
200                    "10000","250", "500", "1000", "10000"), yaxt="n", las=2)
201
```

```
202            axis(2, at=c(0,.2,.4,.6,.8,1),labels=c(0,.2,.4,.6,.8,1),
203             col.axis="black", las=2)
204
205
206            #legend(x="topright", col=2:3, legend=c("detect.MIMIC",
207            #"Factor Analysis"), pch=16, xpd=TRUE)
208
209
210            barplot(cbind(score.list[[i]]$mimic.scores[,2],
211                score.list[[i]]$fa.scores[,2]),
212              main=paste("Graph ",i, " False Positive Rate", sep=""),
213               col=c(2,2,2,2, 3,3,3,3),
214               ylab="Rate", xlab="Number of observations",
215               ylim=c(0,1.2), beside=T, names=c("250", "500", "1000",
216               "10000","250", "500", "1000", "10000"), yaxt="n", las=2)
217
218               axis(2, at=c(0,.2,.4,.6,.8,1),labels=c(0,.2,.4,.6,.8,1),
219                col.axis="black", las=2)
220
221
222               legend(x="topright", col=2:3, legend=c("detect.MIMIC",
223                "Factor Analysis"), pch=16, xpd=TRUE)
224
225            barplot(cbind(score.list[[i]]$mimic.scores[,3],
226                score.list[[i]]$fa.scores[,3]),
227              main=paste("Graph ",i, " True Discovery Rate", sep=""),
228               col=c(2,2,2,2, 3,3,3,3),
229               ylab="Rate", xlab="Number of observations",
230               ylim=c(0,1.2), beside=T, names=c("250", "500", "1000",
231               "10000","250", "500", "1000", "10000"), yaxt="n", las=2)
232
233               axis(2, at=c(0,.2,.4,.6,.8,1),labels=c(0,.2,.4,.6,.8,1),
234                col.axis="black", las=2)
235
236
237         #   legend(x="topright", col=2:3, legend=c("detect.MIMIC",
238         #     "Factor Analysis"), pch=16, xpd=TRUE)
239
240            barplot(cbind(score.list[[i]]$n.null.mimic/500,
241                score.list[[i]]$n.null.fa/500),
242              main=paste("Graph ",i,
243              " Precentage of False\n Latent Cases (out of 500)",
244               sep=""), col=c(2,2,2,2, 3,3,3,3),
245              names=c("250", "500", "1000",
246               "10000","250", "500", "1000", "10000"),
247              ylab="Percentage incorrect",
248              xlab="Number of observations",
249               ylim=c(0,1.2), beside=T, yaxt="n", las=2)
250
251               axis(2, at=c(0,.2,.4,.6,.8,1),labels=c(0,.2,.4,.6,.8,1),
252                col.axis="black", las=2)
253
254         #   legend(x="topright", col=2:3, legend=c("detect.MIMIC",
255         #     "Factor Analysis"), pch=16, xpd=TRUE)
256
257
258         }
259  #       par(mfrow=c(2,2))
260
261  #       for(i in 1:length(score.list)){
262  #           plot(read.dag(paste("sim.graph.", i, ".r.txt", sep="")),
263  #            main=paste("True Graph ", i, sep=""))
264  #
265  #       }
266
267  }
```

100

```r
1
2    # TODO: If neccessary, purge variables after initial pc model is found if they
3    # have both an indegree and outdegree of 0. (for memory optimization
4    # purposes)
5    # TODO: Clean up code.
6    # TODO: Break up long functions into smaller helper functions
7    # TODO: Comment this spaghetti code
8
9    # require(pcalg)
10   # require(igraph)
11   # require(gRim)
12   # require(plotrix)
13
14   find.mimic <- function(data, alpha=.01, indepTest=gaussCItest, pval=.05,
15       print.intermediate=FALSE, high.dim=FALSE){
16
17       orig.names <- names(data)
18
19       names(data) <- paste("X", 1:ncol(data), sep="")
20
21       pc.model <- find.pc.model(data=data, alpha=alpha, indepTest=indepTest)
22
23       deg.model<- graph::degree((pc.model))
24       # If no edges are directed, then return the undirected pc graph.
25
26       if(length(deg.model) < 2 ||
27       length(deg.model) >2){return(pc.model)}
28
29       input.outputs <- find.in.out(pc.model)
30
31       # If no outputs have been found
32       if(is.null(input.outputs$inputs) ||
33           is.null(input.outputs$outputs)){return(pc.model)}
34
35       latent.structure <- finding.latent.structure(input.outputs, pc.model)
36
37       if(high.dim){rm(pc.model)}
38
39       sobers.step <- sobers.criterion(latent.structure, data,
40           input.outputs, pval)
41
42       if(high.dim){rm(latent.structure)}
43
44       last.pc <- final.pc.run(data=data, alpha=alpha, indepTest=indepTest)
45
46       mimic.model.list<-convert.list.to.adj.mat(list.obj=sobers.step,
47           inputs.and.outputs=input.outputs, var.names=names(data))
48
49       n.lat <- ncol(mimic.model.list)
50
51
52       names(mimic.model.list) <- c(orig.names, paste("L",
53        1:(ncol(mimic.model.list)-ncol(data)), sep=""))
54
55
56       mimic.model.graph <- igraph.to.graphNEL(graph.adjacency(as.matrix(mimic.model.list)))
57
58       if(high.dim){rm(mimic.model.list); rm(sobers.step)}
59
60           if(print.intermediate){
61               print("inputs and outputs")
62               print(input.outputs)
63               print("latent.structure - in find.MIMIC")
64               print(latent.structure)
65               print("sobers.step")
66               print(sobers.step)
67               print("mimic.model.list")
```

```r
68                print(mimic.model.list)
69            }
70
71        final.model <- last.step(inputs.outputs=input.outputs, pc.graph=last.pc,
72            mimic.graph=mimic.model.graph)
73
74
75            nodes(final.model) <- c(orig.names, paste("L",
76             1:(n.lat-ncol(data)), sep=""))
77
78
79        if(high.dim){return(list(final.model=final.model))}
80        else{
81            # names(data) <- c(orig.names)
82
83            pre.sober.model <- convert.list.to.adj.mat(list.obj=latent.structure,
84                    inputs.and.outputs=input.outputs, var.names=names(data))
85
86
87            # names(pre.sober.model) <- nodes(final.model)
88
89            pre.sober.model <-
90             igraph.to.graphNEL(graph.adjacency(as.matrix(pre.sober.model)))
91
92            return(list("pc.depth.0"=pc.model, inputs.outputs=input.outputs,
93             latent.structure=latent.structure, pre.sober.model=pre.sober.model,
94            sobers.step=sobers.step, last.pc=last.pc,
95             mimic.model.graph=mimic.model.graph, final.model=final.model))
96        }
97    }
98
99    # Finds PC model. Determines optimal depth via recursion.
100   find.pc.model<-function(data, depth=0, prev.graph=0, indepTest=gaussCItest,
101        alpha=0.01, suffStat=0, n=0, p=0){
102
103        n <- nrow(data)
104        p <- ncol(data)
105
106        ## define sufficient statistics
107        suffStat <- list(C = cor(data), n = n)
108        pc.model<-pc(suffStat=suffStat, indepTest=indepTest, p=p, alpha=alpha,
109                m.max=depth)@graph
110        return(pc.model)
111   }
112
113
114   # Determines inputs/outputs. Takes a GraphNEL object as input. returns names of inputs and outputs
115   find.in.out <- function(graph){
116        indegree.0<- graph::degree(graph)$inDegree==0
117        inputs <- c(which(indegree.0))
118
119        candidate.outputs <- which(!indegree.0)
120
121        adj.matrix <- as(graph, "matrix")
122
123        adj.matrix <- adj.matrix[,candidate.outputs]
124
125        outputs <- apply(adj.matrix, 2,
126            function(pos.output){
127                if(sum(which(pos.output==1)%in%inputs)>0){
128                    return(pos.output)
129                }})
130   if(class(outputs)=="list"){                    102
131        outputs<-remove.null.from.list(outputs)
132        # Extracts the list vector, containing the names of each output.
133        outputs<-names(outputs)
134   }
```

```r
135        else{
136            # Extracts the second list vector, containing the col names of matrix
137            outputs <- dimnames(outputs)[[2]]
138        }
139        inputs <- names(inputs)
140        return(list(inputs=inputs, outputs=outputs))
141    }
142
143    finding.latent.structure <- function(inputs.and.outputs, graph){
144        latent.list <- finding.latents(inputs.and.outputs, graph)
145        latent.list <- remove.null.from.list(latent.list)
146
147            if(length(latent.list)==0){
148                adj.matrix <- as(graph, "matrix")
149                var.names <- names(data.frame(adj.matrix))
150
151                latent.list[['1']]$inputs <-
152                 var.names[as.numeric(inputs.and.outputs$inputs)]
153
154                latent.list[['1']]$outputs <-
155                 var.names[as.numeric(inputs.and.outputs$outputs)]
156
157
158                latent.list[['1']]$latent <- c(1,1)
159                return(latent.list)
160            }
161
162        for(i in 1:(length(latent.list))){
163
164            latent.pair <- smallest.two.subsets(latent.list,
165                    n.inputs=length(inputs.and.outputs$inputs))
166
167            smallest <- latent.pair$smallest
168            nextSmallest <- latent.pair$nextSmallest
169
170            if(smallest==nextSmallest){next()}
171
172                for(j in 1:length(latent.list)) {
173                    if(j!=smallest){
174                        latent.list[[j]]$inputs <-
175                         remove.subset(latent.list[[smallest]]$inputs,
176                                latent.list[[j]]$inputs)
177
178                    }
179
180                }
181
182                latent.list[[nextSmallest]]$latent <-
183                 c(latent.list[[nextSmallest]]$latent,
184                    as.character(smallest), as.character(nextSmallest))
185            }
186        return(latent.list)
187    }
188
189    finding.latents <- function(inputs.and.outputs, graph){
190        adj.matrix <- as(graph, "matrix")
191        var.names <- names(data.frame(adj.matrix))
192
193        inputs  <- inputs.and.outputs$inputs
194        outputs <- inputs.and.outputs$outputs
195
196        input.parents <- adj.matrix[inputs,]
197
198        if(class(input.parents)=="matrix"){
199        #TODO: Redundant. Check if this section is actually called (or is ever called). Data analysis
    version needed changed (had to be fixed).
200            input.parents <- data.frame(input.parents[,
```

103

```r
201                 unique(which(as(input.parents, "matrix")==1, arr.ind=T)[,2])])
202         }
203         else{
204             input.parents <- data.frame(matrix(which(input.parents==1), nrow=1),
205              row.names=inputs)
206             names(input.parents) <- paste("X", input.parents, sep="")
207             input.parents[(input.parents>0)]<-1
208         }
209         latent.list <- construct.latent.list(input.parents,
210                 var.names=var.names)
211
212         return(latent.list)
213     }
214
215     construct.latent.list <- function(input.parents, var.names){
216         latent.list <- c()
217
218         for(i in 1:ncol(input.parents)){
219
220             if(list.exactly.contains(latent.list,
221                 names(col.same(input.parents, column=i)))){next()}
222             else{
223
224                 outputs <- names(input.parents)[col.same(input.parents, column=i)]
225                 latent.list[[i]] <- list(outputs=outputs,
226                     inputs=var.names[as.numeric(get.row.names(input.parents,
227                         outputs))])
228             }
229         }
230         return(latent.list)
231     }
232
233     sobers.criterion <- function(latent.structure, data, inputs.and.outputs,
234         pval=.05){
235         inputs<-c()
236         outputs<-c()
237
238             latents <- get.latents(latent.structure)
239
240             for(i in 1:length(latents)){
241                 if(is.null(latents)){break()}
242                 if(is.null(latent.structure[[i]]$latent)){next()}
243
244                 inputs <- c(get.inputs.via.latents(latent.structure,
245                     latents[[i]]))
246                 outputs <- c(get.outputs.via.latents(latent.structure,
247                         latents[[i]]))
248
249                 inputs <- unique(inputs)
250                 outputs <- unique(outputs)
251
252                 dsep.inputs <- find.dsep(inputs, outputs, data, pval)
253                 if(is.null(dsep.inputs)){return(latent.structure)}
254                 min.set <- 1
255
256                 for(j in 1:length(dsep.inputs)){
257                     if(length(dsep.inputs[[j]]) < length(dsep.inputs[[min.set]])){
258                         min.set<-j
259                     }
260                 }
261                     latent.structure[[as.numeric(latent.structure[[i]]$latent[2])]]$inputs <-
262                  unique(c(dsep.inputs[min.set], latent.structure[[as.numeric(
263                     latent.structure[[i]]$latent[1]]]$inputs))
264                                 latent.structure[[as.numeric(latent.structure[[i]]$latent[1])]]$inputs <-
265                 unique(c(dsep.inputs[min.set], latent.structure[[as.numeric(
266                     latent.structure[[i]]$latent[1])]]$inputs))
```

```r
267
268
269                        if(length(latent.structure[[i]]$latent)==2){
270                            latent.structure[[i]]$latent<-NULL
271                        }
272                        else{
273                            latent.structure[[i]]$latent <- c(latent.structure[[i]]$
274                                latent[3:length(latent.structure[[i]]$latent)])
275
276                        }
277
278            }
279            return(latent.structure)
280  }
281
282  # Last step of the PC run
283  final.pc.run <- function(data, depth=0, prev.graph=NULL,
284        indepTest=gaussCItest, alpha=0.01, suffStat=0, n=0, p=0){
285
286      if(depth==0){
287          n <- nrow(data)
288          p <- ncol(data)
289
290          ## define sufficient statistics
291          suffStat <- list(C = cor(data), n = n)
292      }
293
294      new.graph <- pc(suffStat=suffStat, indepTest=indepTest, p=p, alpha=alpha,
295          m.max=depth+1)@graph
296
297      if(is.null(prev.graph)){prev.graph<-pc(suffStat=suffStat,
298          indepTest=indepTest, p=p, alpha=alpha, m.max=depth)@graph}
299
300      prev.graph <- igraph.from.graphNEL(prev.graph)
301      new.graph <- igraph.from.graphNEL(new.graph)
302
303      if(isTRUE(all.equal(as(as.undirected(prev.graph), "matrix"),
304       as(as.undirected(new.graph), "matrix")))){
305          return(igraph.to.graphNEL(prev.graph))
306      }
307      else{
308          new.graph <- igraph.to.graphNEL(new.graph)
309          return(final.pc.run(data=data, depth=depth+1, prev.graph=new.graph,
310              indepTest=indepTest, alpha=alpha, suffStat=suffStat, n=n, p=p))
311      }
312
313  }
314
315  # Converts Sobers step to adj mat.
316  # TODO: Break into helper functions
317  convert.list.to.adj.mat <- function(list.obj, inputs.and.outputs, var.names){
318      inputs <- inputs.and.outputs$inputs
319      outputs <- inputs.and.outputs$outputs
320      n.variables <- length(var.names)
321      n.latent <- c()
322      n.unconnected.latents <- 0
323
324
325      for(i in 1:length(list.obj)){
326          if(!is.null(list.obj[[i]]$latent)){
327              n.latent<-(c(n.latent, list.obj[[i]]$latent))
328          }
329      }
330
331      for(i in 1:length(list.obj)){
332          if(is.null(list.obj[[i]]$latent) &&
333          isFALSE(as.character(i)%in%n.latent)){
```

105

```r
334                     n.unconnected.latents <- n.unconnected.latents+1
335             }
336         }
337
338     n.latent<-length(unique(n.latent))+n.unconnected.latents
339
340     adj.mat <- matrix(nrow=length(var.names)+n.latent,
341      ncol=length(var.names)+n.latent, data=rep(FALSE,
342         (length(var.names)+n.latent)*(length(var.names)+n.latent)))
343
344     adj.mat <- data.frame(adj.mat)
345
346     names(adj.mat) <- c(var.names, 1:n.latent)
347
348     print((var.names))
349     print(c(row.names(adj.mat), 1:n.latent))
350
351     row.names(adj.mat) <- c(var.names, 1:n.latent)
352
353     # assign latents to their positions
354     for(i in 1:n.latent){
355         if(n.latent>=1){
356             adj.mat[n.variables+i,
357              ] <- c(var.names%in%unlist(list.obj[[i]]$outputs), rep(FALSE,
358                 n.latent))
359
360             adj.mat[, n.variables+i
361              ] <- c(var.names%in%unlist(list.obj[[i]]$inputs), rep(FALSE,
362                 n.latent))
363
364
365             if(!is.null(list.obj[[i]]$latent)){
366                 total.latents <- length(list.obj[[i]]$latent)
367
368
369                 for(j in 1:(ceiling(total.latents/2))){
370                     # As latents are stored in ordered, pairs, ensures that
371                     # the odd position=left, even position=right
372                     left.lat <-list.obj[[i]]$latent[2*(j)-1]
373                     right.lat <- list.obj[[i]]$latent[2*(j)]
374
375                     adj.mat[length(var.names)+as.numeric(left.lat),
376                      length(var.names)+as.numeric(right.lat)]<-TRUE
377                 }
378             }
379         }
380     }
381     # Removes self-causing latents
382     diag(adj.mat)<-FALSE
383     return(adj.mat)
384 }
385
386 # Removes latent-to-output edges based on output of PC depth>0.
387 last.step <- function(inputs.outputs, pc.graph, mimic.graph){
388
389     inputs <- as.numeric(inputs.outputs$inputs)
390     outputs <- as.numeric(inputs.outputs$outputs)
391
392     pc.adj.matrix <- data.frame(as(pc.graph, "matrix"))
393     mimic.adj.matrix <- data.frame(as(mimic.graph, "matrix"))
394
395     n.vars <- ncol(pc.adj.matrix)
396
397     n.latents <- ncol(mimic.adj.matrix)-n.vars
398
399     false.outputs <- names(which(apply(pc.adj.matrix[inputs,
400      outputs], 2, sum)==0)))
```

106

```r
401
402        names(mimic.adj.matrix) <- c(names(pc.adj.matrix), paste("L",
403         1:n.latents, sep=""))
404
405
406     if(!is.null(false.outputs) && length(false.outputs)>0){
407            for(j in false.outputs){
408
409                false.connected.latents<-(which(mimic.adj.matrix[
410                 j]==1))
411
412                mimic.adj.matrix[false.connected.latents,
413                 which(names(mimic.adj.matrix)%in%j)] <- 0
414
415            }
416
417        for(i in false.outputs){
418            false.connected.latents<-(which(mimic.adj.matrix[
419             i]==1))
420            mimic.adj.matrix[(1:n.vars),
421             which(names(mimic.adj.matrix)%in%i)] <-
422             pc.adj.matrix[,which(names(mimic.adj.matrix)%in%i)]
423
424    mimic.adj.matrix[which(names(mimic.adj.matrix)%in%i),
425    -(false.connected.latents)] <-
426     pc.adj.matrix[which(names(mimic.adj.matrix)%in%i),]
427
428            }
429
430        }
431    return(igraph.to.graphNEL(graph.adjacency(as.matrix(mimic.adj.matrix))))
432 }
433
434
435
436 # Helper Methods
437
438 # Returns the index of columns that are all identical.
439 col.same <- function(mat, column){
440     return(which(colSums(mat[,column]==mat)==nrow(mat)))
441 }
442
443 list.exactly.contains <- function(list.object, search.term){
444     return(isTRUE(sum(unlist(lapply(list.object,
445        function(item){(search.term %in% item$outputs)})))>0))
446 }
447
448 get.row.names <- function(mat, col.names){
449     return(row.names(mat)[unique(which(mat[c(col.names)]==1, arr.ind=T)[,1])])
450 }
451
452 remove.null.from.list <- function(list.object){
453     purged.list <- list.object[-which(sapply(list.object,
454            is.null),arr.ind=TRUE)]
455
456     if(length(purged.list)==0){return(list.object)}
457     else{return(purged.list)}
458 }
459
460 smallest.two.subsets <- function(latent.list, n.inputs){
461     smallest <- NULL
462     nextSmallest <- NULL
463                                    107
464     for(i in 1:(length(latent.list))){
465         for(j in 1:length(latent.list)){
466
467            if(i==j){next()}
```

```r
468
469              if(is.subset(latent.list[[i]]$input, latent.list[[j]]$input)){
470
471                  # Ensures that subsets are being compared for smallest/n.small
472                  if(is.null(smallest)){smallest<-i; nextSmallest<-j}
473
474                  if(isTRUE(length(latent.list[[i]]$inputs) <=
475                   length(latent.list[[smallest]]$inputs))){smallest<-i}
476                  else if(isTRUE(length(latent.list[[i]]$inputs) <=
477                   length(latent.list[[nextSmallest]]$inputs))){nextSmallest<-i}
478
479              }
480              else if (is.subset(latent.list[[j]]$input,
481                   latent.list[[i]]$input)){
482
483                  # Ensures that subsets are being compared for smallest/n.small
484                  if(is.null(smallest)){smallest<-j; nextSmallest<-i}
485
486                  if(isTRUE(length(latent.list[[j]]$inputs) <=
487                   length(latent.list[[smallest]]$inputs))){smallest<-j}
488                  else if(isTRUE(length(latent.list[[j]]$inputs) <=
489                   length(latent.list[[nextSmallest]]$inputs))){nextSmallest<-j}
490              }
491          }
492      }
493
494      if(is.null(smallest)){return(list(smallest=1, nextSmallest=1))}
495
496      return(list(smallest=smallest, nextSmallest=nextSmallest))
497  }
498
499  there.is.a.subset <- function(latent.list, n.inputs){
500      size.list <- smallest.two.subsets(latent.list,
501          n.inputs=length(inputs.and.outputs$inputs))
502          return(size.list[[1]]==size.list[[2]])
503
504  }
505
506
507  remove.subset <- function(small.set, larger.set){
508      return(larger.set[!larger.set%in%small.set])
509  }
510
511  is.subset <- function(set.1, set.2){
512
513      if(length(set.1)>length(set.2)){return(FALSE)}
514
515      joint.membership <- c(0)
516      for(element in set.1){
517          joint.membership <- c(joint.membership + sum(set.2==element))
518      }
519
520      if(joint.membership==length(set.1)){
521          return(TRUE)
522      }
523      else{
524          return(FALSE)
525      }
526
527  }
528
529  get.latents <- function(latent.structure){
530      latents <- c()
531      for(i in 1:length(latent.structure)){
532          latents[[i]] <- latent.structure[[i]]$latent
533      }
534      return(latents)
```

```r
535    }
536
537
538    get.inputs.via.latents <- function(latent.structure, latents){
539        latents <- as.numeric(latents)
540        inputs <- c()
541
542        left.latent <- latents[1]
543        right.latent <- latents[2]
544
545        inputs <- c(latent.structure[[left.latent]]$inputs,
546            latent.structure[[right.latent]]$inputs)
547
548        return(inputs)
549    }
550
551    get.outputs.via.latents <- function(latent.structure, latents){
552        latents <- as.numeric(latents)
553        outputs <- c()
554
555        left.latent <- latents[1]
556        right.latent <- latents[2]
557
558        outputs <- c(latent.structure[[left.latent]]$outputs,
559            latent.structure[[right.latent]]$outputs)
560
561        return(outputs)
562    }
563
564    find.dsep <- function(inputs, outputs, data, pval=.05){
565            variations.list <- c()
566            n.inputs <- length(inputs)
567
568            for(i in 1:n.inputs){
569                variations.list[[i]] <- list(combn(inputs, m=i))
570            }
571
572            condi.sets<-variations.list
573
574            for(i in 1:length(condi.sets)){
575                dsep.set<-get.cond.combo(data, outputs, condi.sets[i], pval)
576                if(!is.null(dsep.set)){
577                    # print(dsep.set)
578                    return(dsep.set)
579                }
580            }
581        return()
582    }
583
584    get.cond.combo <- function(data, outputs, input, pval){
585
586        input <- destroy.list(input)
587
588        # Handles single input cases.
589        if(is.null(dim(input))){
590
591            for(i in 1:length(input)){
592                if(!ciTest(data, set=c(outputs[c(1, length(outputs))],
593                  (input[i])))$p.val<=pval){return(input[i])}
594            }
595        }
596        else{
597            for(i in 1:ncol(input)){
598                if(!ciTest(data, set=c(outputs[c(1, length(outputs))],
599                  (input[,i])))$p.val<=pval){return(input[,i])}
600            }
601        return(c())
```

109

```
602        }
603    }
604
605    destroy.list <- function(list.obj){
606        if(class(list.obj)=="list"&&
607            listDepth(list.obj)>1){
608
609            destroy.list(unlist(list.obj, recursive=F))
610        }
611        else{return(list.obj[[1]])}
612    }
613
614    isFALSE <- function(truth.vector){return(!isTRUE(truth.vector))}
615
```

```
1   source("new\ code/high_dim_indep.R")
2   #Running steps:
3   library(plotrix)
4
5   #Data formating
6
7   # Reading in output variables.
8   ov_tumor_expression.df <- read.csv("OV_tumor_expression_matrix_fold.csv")
9   ov_tumor_expression.df <-
10   apply(ov_tumor_expression.df[,2:ncol(ov_tumor_expression.df)], 2,
11   as.numeric)
12
13  #Reading in input variables.
14  ov_tumor_sga.df <- read.csv("OV_tumor_SGA_matrix.csv")
15  ov_tumor_sga.df <- apply(ov_tumor_sga.df[,2:ncol(ov_tumor_sga.df)], 2,
16   as.numeric)
17
18
19  # Converting both datasets into factor variables.
20  ov_tumor_expression.df<-apply( ov_tumor_expression.df, 2, function(data){factor(data, levels=c
    (-1,0,1))})
21  ov_tumor_sga.df<-apply( ov_tumor_sga.df, 2, function(data){factor(data, levels=c(0,1))})
22
23  #converting from matrix to data.frame
24  ov_tumor_expression.df <- data.frame(ov_tumor_expression.df)
25  ov_tumor_sga.df <- data.frame(ov_tumor_sga.df)
26
27  #Combining datasets into a single dataset.
28  ov_tumor.df <- data.frame(ov_tumor_sga.df, ov_tumor_expression.df)
29
30  library(glmnet)
31
32  #converts to needed format for glmnet
33  ov_tumor.matrix <- data.matrix(ov_tumor.df)
34
35  is.constant <- function(vect){
36      return(sum(vect[1]==vect)==length(vect))
37  }
38
39  inputs <- names(ov_tumor_sga.df)
40  inputs <- 1:length(inputs)
41  outputs <- names(ov_tumor_expression.df)
42  outputs <- (length(inputs)+1):(ncol(ov_tumor.df))
43
44  #Gets list of non-constant variables
45  exclude.constants<-!apply(ov_tumor.df, 2, is.constant)
46  exclude.constants<-which(exclude.constants)
47
48  inputs.cleaned <- which(exclude.constants<=max(inputs))
49  outputs.cleaned <- which(exclude.constants>max(inputs))
50
51  library(doParallel)
52  registerDoParallel(8)
53
54  #removing inputs
55  cv.lasso.results<-list()
56  for(i in (outputs.cleaned)){
57      cv.lasso.results[[i]]<-cv.glmnet(alpha=1, x=ov_tumor.matrix[,inputs.cleaned], y=ov_tumor.matrix
    [,i], standardize=F, parallel=TRUE)
58
59      cv.lasso.results[[i]]<-coef(cv.lasso.results[[i]], s=cv.lasso.results[[i]]$lambda.min)
60
61      # The [[1]] part of the code selects the correct part of the Dimnames object. See str
    (cv.lasso.results[[i]]) for why.
62      cv.lasso.results[[i]]<- cv.lasso.results[[i]]@Dimnames[[1]][which(cv.lasso.results[[i]]>0)]
63  }
64
```

```r
65
66
67   # removing outputs
68   cv.lasso.results.inputs<-list()
69   for(i in (inputs.cleaned)){
70       cv.lasso.results.inputs[[i]]<-cv.glmnet(alpha=1, x=ov_tumor.matrix[,outputs.cleaned],
     y=ov_tumor.matrix[,i], standardize=F, parallel=TRUE)
71
72       cv.lasso.results.inputs[[i]]<-coef(cv.lasso.results.inputs[[i]], s=cv.lasso.results.inputs[[i]]
     $lambda.min)
73
74       # The [[1]] part of the code selects the correct part of the Dimnames object. See str
     (cv.lasso.results[[i]]) for why.
75       cv.lasso.results.inputs[[i]]<- cv.lasso.results.inputs[[i]]@Dimnames[[1]][which
     (cv.lasso.results.inputs[[i]]>0)]
76
77   }
78
79   # Initial attempt
80   #reduced.dataset<-ov_tumor.df[,(names(ov_tumor.df[, inputs])%in%(unique(unlist(cv.lasso.results))))]
81   #inputs<-1:length(ov_tumor.df[,(names(ov_tumor.df[, inputs])%in%(unique(unlist
     (cv.lasso.results)))))])
82
83   #reduced.dataset<-data.frame(reduced.dataset, ov_tumor.df[,(names(ov_tumor.df[, outputs])%in%(unique
     (unlist(cv.lasso.results.inputs)))))])
84   #outputs<-(length(ov_tumor.df[,(names(ov_tumor.df[, inputs])%in%(unique(unlist
     (cv.lasso.results)))))])+1):ncol(reduced.dataset)
85
86   # Need even more:
87   # 95% quantile: 42, 111, drops to 7085 variables.
88   # 99% quantile: 69, 141, drops to 4369 variables. verified.
89   temp<-cv.lasso.results[unlist(lapply(cv.lasso.results, length)>69)]
90   temp2<-cv.lasso.results.inputs[unlist(lapply(cv.lasso.results.inputs, length)>141)]
91
92   #Need to fix, as the renaning using .1 leads to too many variables being missed.
93   reduced.dataset.final<-ov_tumor.df[,(names(ov_tumor.df)%in%(unique(unlist(temp))))]
94   inputs.final<-1:length(ov_tumor.df[,(names(ov_tumor.df)%in%(unique(unlist(temp))))])
95
96   reduced.dataset.final<-data.frame(reduced.dataset.final, ov_tumor.df[,(names(ov_tumor.df)%in%(unique
     (unlist(temp2))))])
97   outputs.final<-(length(ov_tumor.df[,(names(ov_tumor.df)%in%(unique(unlist(temp))))])+1):ncol
     (reduced.dataset.final)
98
99
100  # Checking to see that have correctly subseted dataset. Note that 2 is subtracted as the
     (intercept) carries over.
101  #length(unique(unlist(temp2)))+length(unique(unlist(temp)))-2
102  #dim(reduced.dataset.final)
103
104
105  # This dataset has 4367 variables.
106  write.csv(reduced.dataset.final, "reduced.datasetv2.csv", row.names=F)
107
108  orig.names <-names(reduced.dataset.final)
109  names(reduced.dataset.final) <-paste("X", 1:ncol(reduced.dataset.final), sep="")
110
111  initial.adj.mat <- alternative.to.pc(inputs.final, outputs.final, dataset.df=reduced.dataset.final,
     p.val=.000001)
112
113  init.latent.struct <- finding.latent.structure(inputs.and.outputs=list(inputs=inputs.final,
     outputs=outputs.final), adj.matrix=initial.adj.mat)
114
115  sobers.step <- sobers.criterion(init.latent.struct, reduced.dataset.final,
116          inputs.and.outputs=list(inputs=inputs.final, outputs=outputs.final), pval=.000001,
     categorical=TRUE)
117
118  #Converts list type object to desired data.frame version of adj.matrix. Data frame is used so that
```

```r
        variables can easily be renamed.
119   mimic.model.list<-convert.list.to.adj.mat(list.obj=sobers.step,
120           inputs.and.outputs=list(inputs=inputs.final, outputs=outputs.final), var.names=names
      (reduced.dataset.final))
121
122
123   n.lat <- ncol(mimic.model.list)
124
125   #Renames variables to preserve original variable names
126   names(mimic.model.list) <- c(orig.names, paste("L",
127       1:(ncol(mimic.model.list)-ncol(data)), sep=""))
128
129   library(igraph)
130   #Coverts dataframe to graphNEL object.
131   mimic.model.graph <- igraph.to.graphNEL(graph.adjacency(as.matrix(mimic.model.list)))
132
133
134   #Adds back in the variable names. 1:139 names the latents.
135   nodes(mimic.model.graph) <- c(orig.names, 1:139)
136
137   # Number of nodes with an indegree>0
138   summary(graph::degree(mimic.model.graph)$inDegree>0)
139
140   # Number of nodes with an outdegree>0
141   summary(graph::degree(mimic.model.graph)$outDegree>0)
142
143   # Number of nodes with a total degree>0
144   length(names(which((graph::degree(mimic.model.graph)$outDegree +graph::degree(mimic.model.graph)
      $inDegree)>0)))
145
146   #Extracting a subgraph only containing nodes connected to other nodes.
147   reduced.final.graph<-igraph.from.graphNEL(subGraph(snodes=names(which((graph::degree
      (mimic.model.graph)$outDegree +graph::degree(mimic.model.graph)$inDegree)>0)),
      graph=mimic.model.graph))
148
149
150   latent.nodes<-which(V(reduced.final.graph)$name%in%1:n.lat)
151   input.nodes <- which(V(reduced.final.graph)$name%in%(names(which(degree(reduced.final.graph,
      mode="in")==0))))
152   output.nodes <- which(V(reduced.final.graph)$name%in%(names(which(degree(reduced.final.graph,
      mode="out")==0))))
153
154   #Assigning color based on whether node is input, latent, or output node. 1, 2, and 3 are colors.
155   cluster.mem.color <- ifelse(1:length(V(reduced.final.graph)$name)%in%latent.nodes, 1,
156               ifelse(1:length(V(reduced.final.graph)$name)%in%output.nodes, 3,
157                   ifelse(1:length(V(reduced.final.graph)$name)%in%input.nodes, 2, 0)))
158
159   # Full protein signalling network.
160   pdf("signal_network.pdf")
161   plot.igraph(reduced.final.graph, vertex.label=NA, edge.width=1, edge.arrow.size=0, vertex.size=3,
      vertex.color=cluster.mem.color, vertex.frame.color=NA, layout=layout.fruchterman.reingold
      (reduced.final.graph),
162       main="Protein Signalling Network")
163   legend("topright", pch=16, col=c(1,2,3), legend=c("Latent","Input","Output"))
164   dev.off()
165
166
167
168   # Extracts a subgraph built from a given starting node.
169   plot.interesting.subgraph <- function(reduced.final.graph, starting.nodes=c("75"), ...){
170
171       reduced.subgraph <- induced.subgraph(reduced.final.graph, v=get.subgraph(reduced.final.graph,
      starting.nodes=starting.nodes))
172
173
174       latent.nodes<-which(V(reduced.subgraph)$name%in%1:(length(V(reduced.final.graph))))
175
```

113

```R
176        input.nodes <- which(V(reduced.subgraph)$name%in%(names(which(degree(reduced.subgraph,
    mode="in")==0))))
177
178        output.nodes <- which(V(reduced.subgraph)$name%in%(names(which(degree(reduced.subgraph,
    mode="out")==0))))
179
180        cluster.mem.color <- ifelse(1:length(V(reduced.subgraph)$name)%in%latent.nodes, 1,
181                      ifelse((1:length(V(reduced.subgraph)$name)%in%output.nodes), 3,
182                          ifelse((1:length(V(reduced.subgraph)$name)%in%input.nodes), 2, 0)))
183
184        plot.igraph(reduced.subgraph, layout=layout.fruchterman.reingold(reduced.subgraph,
    area=40*vcount(reduced.subgraph)^2), vertex.label.color=cluster.mem.color, ...)
185   }
186
187
188
189   # Extracts a the maximal subgraph were all of the nodes are connected to one another, given a
    starting node to search from.
190   get.subgraph <- function(orig.graph, starting.nodes=c(), searched.nodes=c()){
191        adj.nodes<-c()
192        for(i in 1:length(starting.nodes)){
193
194            new.adj <- getElement(get.adjlist(reduced.final.graph), starting.nodes[i])
195            adj.indices <- new.adj
196
197            adj.nodes <- unique(c(adj.nodes, V(orig.graph)$name[as.numeric(adj.indices)]))
198        }
199
200        if(sum(adj.nodes %in% searched.nodes)==length(adj.nodes)){
201            return(unique(searched.nodes))
202        }
203        else{
204            return(unique(get.subgraph(orig.graph, starting.node=adj.nodes, searched.nodes=unique(c
    (adj.nodes, searched.nodes)))))
205        }
206   }
207
208   # Plots an interesting subgraph.
209   pdf("subgraph_final.pdf")
210   plot.interesting.subgraph(reduced.final.graph, edge.width=1, vertex.label=c(), edge.arrow.size=0.3,
    vertex.label.cex=.5, vertex.color=0,
211   vertex.size=7, vertex.frame.color=NA, main="Protein Signalling Network", starting.nodes="76")
212   legend("topright", pch=16, col=c(1,2,3), legend=c("Latent","Input","Output"))
213   dev.off()
214
215
216   #Plot all subgraphs with 4 or more variables.
217   pdf("all_subgraph_final.pdf")
218   already.tested<-c()
219        for(i in 1:length(latent.nodes)){
220        temp.subgraph <- get.subgraph(reduced.final.graph, starting.nodes=i)
221        if(length(temp.subgraph)<4 || sum(temp.subgraph%in%unlist(already.tested))==length
    (temp.subgraph)){next()}
222            else{
223            plot.interesting.subgraph(reduced.final.graph, edge.width=1, vertex.label=c(),
    edge.arrow.size=0.3, vertex.label.cex=.5,
224             vertex.color=0,vertex.size=7, vertex.frame.color=NA, main="Protein Signalling Network",
    starting.nodes=(i))
225            legend("topright", pch=16, col=c(1,2,3), legend=c("Latent","Input","Output"))
226
227            already.tested[[i]]<-temp.subgraph
228
229        }
230   }
231   dev.off()
232
233   #EDA
```

114

```r
234
235    pdf("eda_final.pdf")
236    # N.abnormalities (numeric vector) per subject
237    freq.of.abnormalities <- apply(ov_tumor_sga.df==1, 1, sum)
238
239    barplot(sort(freq.of.abnormalities), xlab="Subject", ylab="Number of abnormalities", main="Number
       of Mutations\n for each Subject",
240     ylim=c(0,max(freq.of.abnormalities)+100))
241
242
243    # N.subjects (numeric vector) per mutation
244    freq.of.subject <- apply(ov_tumor_sga.df==1, 2, sum)
245
246
247    # Bar Plot of frequency of abnormality.
248    barplot(sort(freq.of.subject), xlab="Gene", ylab="Number of subjects with a particular genetic
       mutation", main="Mutation Frequency for each Gene", ylim=c(0,
249        max(freq.of.subject)+100), names.arg=NA)
250
251    dev.off()
```

```r
1    require(gRim)
2
3    #require(pcalg)
4    #Takes vector containg index of inputs in dataset.df (same for outputs).
5    alternative.to.pc <- function(inputs, outputs, dataset.df, p.val){
6        adj.mat<-matrix(nrow=ncol(dataset.df), ncol=ncol(dataset.df), 0)
7
8        for(i in min(inputs):max(inputs)){
9        #get number of levels for factor
10           n.levels.input <- length(levels((dataset.df[,i])))
11
12           for(j in min(outputs):max(outputs)){
13               #get number of levels for factor
14               n.levels.output <- length(levels((dataset.df[,j])))
15
16           #If one of the variables is a constant, skip it.
17           if(n.levels.output>1 && n.levels.input>1){
18               # If variables i and j are dependent, write a one in the ith column and jth row of the
     adj. matrix.
19               adj.mat[i, j] <- (fisher.test(x=dataset.df[,i], y=dataset.df[,j])$p.value<=p.val)
20           }
21
22       }
23
24       }
25       return(adj.mat)
26   }
27
28   finding.latent.structure <- function(inputs.and.outputs, adj.matrix){
29
30       var.names <- paste("X", 1:ncol(adj.matrix), sep="")
31
32       latent.list <- finding.latents(inputs.and.outputs, adj.matrix, var.names)
33       latent.list <- remove.null.from.list(latent.list)
34
35           if(length(latent.list)==0){
36               #adj.matrix <- as(graph, "matrix")
37               #var.names <- names(data.frame(adj.matrix))
38
39               latent.list[['1']]$inputs <-
40                var.names[as.numeric(inputs.and.outputs$inputs)]
41
42               latent.list[['1']]$outputs <-
43                var.names[as.numeric(inputs.and.outputs$outputs)]
44
45
46               latent.list[['1']]$latent <- c(1,1)
47               return(latent.list)
48           }
49
50       for(i in 1:(length(latent.list))){
51
52           latent.pair <- smallest.two.subsets(latent.list,
53               n.inputs=length(inputs.and.outputs$inputs))
54
55           smallest <- latent.pair$smallest
56           nextSmallest <- latent.pair$nextSmallest
57
58
59           if(smallest==nextSmallest){next()}
60           else{
61               for(j in 1:length(latent.list)) {
62                   if(j!=smallest  && length(latent.list[[j]]$inputs)>1){
63
64                       latent.list[[j]]$inputs <-
65                        remove.subset(latent.list[[smallest]]$inputs,
66                            latent.list[[j]]$inputs)
```

```R
67
68                            }
69
70                    }
71                }
72
73                latent.list[[nextSmallest]]$latent <-
74                 c(latent.list[[nextSmallest]]$latent,
75                     as.character(smallest), as.character(nextSmallest))
76            }
77        return(latent.list)
78    }
79
80     finding.latents <- function(inputs.and.outputs, adj.matrix, var.names){
81        #adj.matrix <- as(graph, "matrix")
82        #var.names <- names(data.frame(adj.matrix))
83
84
85        inputs   <- inputs.and.outputs$inputs
86        outputs <- inputs.and.outputs$outputs
87
88        #print(inputs)
89
90        # Note: This is a line that is adding alot of memory usage.
91        input.parents <- adj.matrix[inputs,]
92
93        if(class(input.parents)=="matrix"){
94
95            output.names <- paste("X", which(apply(input.parents,2, function(column){sum(column)>0})),
    sep="")
96            input.parents <- data.frame(input.parents[,
97                unique(which(as(input.parents, "matrix")==1, arr.ind=T)[,2])])
98            names(input.parents) <- output.names
99        }
100       else{        #  1 in column => column var. effect of row var where 1 occurs.
101           input.parents <- data.frame(matrix(which(input.parents==1), nrow=1),
102            row.names=inputs)
103           names(input.parents) <- var.names
104           input.parents[(input.parents>0)]<-1
105       }
106
107       latent.list <- construct.latent.list(input.parents,
108              var.names=var.names)
109
110       return(latent.list)
111   }
112
113   construct.latent.list <- function(input.parents, var.names){
114       latent.list <- c()
115
116       for(i in 1:ncol(input.parents)){
117
118           if(list.exactly.contains(latent.list,
119                names(col.same(input.parents, column=i)))){next()}
120           else{
121
122               outputs <- names(input.parents)[col.same(input.parents, column=i)]
123               latent.list[[i]] <- list(outputs=outputs,
124                     inputs=var.names[as.numeric(get.row.names(input.parents,
125                        outputs))])
126           }
127       }
128       return(latent.list)                         117
129   }
130
131   sobers.criterion <- function(latent.structure, data, inputs.and.outputs,
132        pval=.05, categorical=TRUE){
```

```r
133        inputs<-c()
134        outputs<-c()
135
136            latents <- get.latents(latent.structure)
137
138            for(i in 1:length(latents)){
139                if(is.null(latents)){break()}
140                if(is.null(latent.structure[[i]]$latent)){next()}
141
142                inputs <- c(get.inputs.via.latents(latent.structure,
143                    latents[[i]]))
144                outputs <- c(get.outputs.via.latents(latent.structure,
145                    latents[[i]]))
146
147                inputs <- unique(inputs)
148                outputs <- unique(outputs)
149
150                dsep.inputs <- find.dsep(inputs, outputs, data, pval, categorical)
151                if(is.null(dsep.inputs)){return(latent.structure)}
152                min.set <- 1
153
154                for(j in 1:length(dsep.inputs)){
155                    if(length(dsep.inputs[[j]]) < length(dsep.inputs[[min.set]])){
156                        min.set<-j
157                    }
158                }
159                latent.structure[[as.numeric(latent.structure[[i]]$latent[2])]]$inputs <-
160              unique(c(dsep.inputs[min.set], latent.structure[[as.numeric(
161                latent.structure[[i]]$latent[2])]]$inputs))
162
163                latent.structure[[as.numeric(latent.structure[[i]]$latent[1])]]$inputs <-
164              unique(c(dsep.inputs[min.set], latent.structure[[as.numeric(
165                latent.structure[[i]]$latent[1])]]$inputs))
166
167
168                if(length(latent.structure[[i]]$latent)==2){
169                    latent.structure[[i]]$latent<-NULL
170                }
171                else{
172                    latent.structure[[i]]$latent <- c(latent.structure[[i]]$
173                        latent[3:length(latent.structure[[i]]$latent)])
174
175                }
176
177            }
178            return(latent.structure)
179    }
180
181
182    # Converts Sobers step to adj mat.
183    # TODO: Break into helper functions.
184    # TODO: Fix memory stuff in this.
185    convert.list.to.adj.mat <- function(list.obj, inputs.and.outputs, var.names){
186        inputs <- inputs.and.outputs$inputs
187        outputs <- inputs.and.outputs$outputs
188        n.variables <- length(var.names)
189        n.latent <- c()
190        n.unconnected.latents <- 0
191
192        for(i in 1:length(list.obj)){
193            if(!is.null(list.obj[[i]]$latent)){
194                n.latent<-(c(n.latent, list.obj[[i]]$latent))
195            }
196        }
197
198        for(i in 1:length(list.obj)){
199            if(is.null(list.obj[[i]]$latent) &&
```

118

```r
200                 isFALSE(as.character(i)%in%n.latent)){
201                     n.unconnected.latents <- n.unconnected.latents+1
202                 }
203         }
204
205         n.latent<-length(unique(n.latent))+n.unconnected.latents
206
207         adj.mat <- matrix(nrow=length(var.names)+n.latent,
208          ncol=length(var.names)+n.latent, data=rep(FALSE,
209             (length(var.names)+n.latent)*(length(var.names)+n.latent)))
210
211         adj.mat <- data.frame(adj.mat)
212
213         names(adj.mat) <- c(var.names, 1:n.latent)
214
215         print((var.names))
216         print(c(row.names(adj.mat), 1:n.latent))
217
218         row.names(adj.mat) <- c(var.names, 1:n.latent)
219
220         # assign latents to their positions
221         for(i in 1:n.latent){
222             if(n.latent>=1){
223                 adj.mat[n.variables+i,
224                  ] <- c(var.names%in%unlist(list.obj[[i]]$outputs), rep(FALSE,
225                     n.latent))
226
227                 adj.mat[, n.variables+i
228                  ] <- c(var.names%in%unlist(list.obj[[i]]$inputs), rep(FALSE,
229                     n.latent))
230
231
232                 if(!is.null(list.obj[[i]]$latent)){
233                     total.latents <- length(list.obj[[i]]$latent)
234
235
236                     for(j in 1:(ceiling(total.latents/2))){
237                         # As latents are stored in ordered, pairs, ensures that
238                         # the odd position=left, even position=right
239                         left.lat <-list.obj[[i]]$latent[2*(j)-1]
240                         right.lat <- list.obj[[i]]$latent[2*(j)]
241
242                         adj.mat[length(var.names)+as.numeric(left.lat),
243                          length(var.names)+as.numeric(right.lat)]<-TRUE
244                     }
245                 }
246             }
247         }
248         # Removes self-causing latents
249         diag(adj.mat)<-FALSE
250         return(adj.mat)
251 }
252
253 # Helper Methods
254
255 # Returns the index of columns that are all identical.
256 col.same <- function(mat, column){
257     return(which(colSums(mat[,column]==mat)==nrow(mat)))
258 }
259
260 list.exactly.contains <- function(list.object, search.term){
261     return(isTRUE(sum(unlist(lapply(list.object,
262         function(item){(search.term %in% item$outputs)})))>0))
263 }
264
265 get.row.names <- function(mat, col.names){
266     return(row.names(mat)[unique(which(mat[c(col.names)]==1, arr.ind=T)[,1])])
```

```R
267   }
268
269   remove.null.from.list <- function(list.object){
270       purged.list <- list.object[-which(sapply(list.object,
271               is.null),arr.ind=TRUE)]
272
273       if(length(purged.list)==0){return(list.object)}
274       else{return(purged.list)}
275   }
276
277   smallest.two.subsets <- function(latent.list, n.inputs){
278       smallest <- NULL
279       nextSmallest <- NULL
280
281       for(i in 1:(length(latent.list))){
282           for(j in 1:length(latent.list)){
283
284               if(i==j){next()}
285
286               if(is.subset(latent.list[[i]]$input, latent.list[[j]]$input)){
287
288                   # Ensures that subsets are being compared for smallest/n.small
289                   if(is.null(smallest)){smallest<-i; nextSmallest<-j}
290
291                   if(isTRUE(length(latent.list[[i]]$inputs) <=
292                    length(latent.list[[smallest]]$inputs))){smallest<-i}
293                   else if(isTRUE(length(latent.list[[i]]$inputs) <=
294                    length(latent.list[[nextSmallest]]$inputs))){nextSmallest<-i}
295
296               }
297               else if (is.subset(latent.list[[j]]$input,
298                   latent.list[[i]]$input)){
299
300                   # Ensures that subsets are being compared for smallest/n.small
301                   if(is.null(smallest)){smallest<-j; nextSmallest<-i}
302
303                   if(isTRUE(length(latent.list[[j]]$inputs) <=
304                    length(latent.list[[smallest]]$inputs))){smallest<-j}
305                   else if(isTRUE(length(latent.list[[j]]$inputs) <=
306                    length(latent.list[[nextSmallest]]$inputs))){nextSmallest<-j}
307               }
308           }
309       }
310
311       if(is.null(smallest)){return(list(smallest=1, nextSmallest=1))}
312
313       return(list(smallest=smallest, nextSmallest=nextSmallest))
314   }
315
316   there.is.a.subset <- function(latent.list, n.inputs){
317       size.list <- smallest.two.subsets(latent.list,
318           n.inputs=length(inputs.and.outputs$inputs))
319           return(size.list[[1]]==size.list[[2]])
320
321   }
322
323
324   remove.subset <- function(small.set, larger.set){
325       return(larger.set[!larger.set%in%small.set])
326   }
327
328   is.subset <- function(set.1, set.2){
329                                       120
330       if(length(set.1)>length(set.2)){return(FALSE)}
331
332       joint.membership <- c(0)
333       for(element in set.1){
```

```r
334           joint.membership <- c(joint.membership + sum(set.2==element))
335       }
336
337       if(joint.membership==length(set.1)){
338           return(TRUE)
339       }
340       else{
341           return(FALSE)
342       }
343
344  }
345
346  get.latents <- function(latent.structure){
347       latents <- c()
348       for(i in 1:length(latent.structure)){
349           latents[[i]] <- latent.structure[[i]]$latent
350       }
351       return(latents)
352  }
353
354
355  get.inputs.via.latents <- function(latent.structure, latents){
356       latents <- as.numeric(latents)
357       inputs <- c()
358
359       left.latent <- latents[1]
360       right.latent <- latents[2]
361
362       inputs <- c(latent.structure[[left.latent]]$inputs,
363           latent.structure[[right.latent]]$inputs)
364
365       return(inputs)
366  }
367
368  get.outputs.via.latents <- function(latent.structure, latents){
369       latents <- as.numeric(latents)
370       outputs <- c()
371
372       left.latent <- latents[1]
373       right.latent <- latents[2]
374
375       outputs <- c(latent.structure[[left.latent]]$outputs,
376           latent.structure[[right.latent]]$outputs)
377
378       return(outputs)
379  }
380
381  find.dsep <- function(inputs, outputs, data, pval=.05, categorical=FALSE){
382           variations.list <- c()
383           n.inputs <- length(inputs)
384
385           for(i in 1:n.inputs){
386               variations.list[[i]] <- list(combn(inputs, m=i))
387           }
388
389           condi.sets<-variations.list
390
391           for(i in 1:length(condi.sets)){
392               dsep.set<-get.cond.combo(data, outputs, condi.sets[i], pval, categorical)
393               if(!is.null(dsep.set)){
394                   # print(dsep.set)
395                   return(dsep.set)
396               }
397           }
398       return()
399  }
400
```

121

```r
401    #TODO comment this fuction. The nested conditionals/flow control especially.
402    get.cond.combo <- function(data, outputs, input, pval, categorical=FALSE){
403
404        input <- destroy.list(input)
405
406        # Handles single input cases.
407        if(is.null(dim(input))){
408
409            for(i in 1:length(input)){
410
411                if(isTRUE(categorical)){
412                    if(!ciTest(data, set=c(outputs[c(1, length(outputs))],
413                    (input[i])))$p.val<=pval){return(input[i])}
414                }
415                else{
416                    if(!ciTest(data, set=c(outputs[c(1, length(outputs))],
417                    (input[i])))$p.val<=pval){return(input[i])}
418                }
419            }
420        }
421        else{
422            for(i in 1:ncol(input)){
423                if(isTRUE(categorical)){
424
425                }
426                else{
427                if(!ciTest(data, set=c(outputs[c(1, length(outputs))],
428                 (input[,i])))$p.val<=pval){return(input[,i])}
429                }
430            }
431        return(c())
432        }
433    }
434
435    destroy.list <- function(list.obj){
436        if(class(list.obj)=="list"&&
437            listDepth(list.obj)>1){
438
439            destroy.list(unlist(list.obj, recursive=F))
440        }
441        else{return(list.obj[[1]])}
442    }
443
444    isFALSE <- function(truth.vector){return(!isTRUE(truth.vector))}
```

122

# Bibliography

Brodie, A. (2014). Identifying endogenous latent causal structure under linearity and sparsity assumptions.

Bühn, A. and Schneider, F. G. (2008). Mimic models, cointegration and error correction: an application to the french shadow economy. Technical report, CESifo working paper.

Craig, E. and Hoskin, M. (1992). Hegel and the seven planets. *Journal for the History of Astronomy*, 23(3).

DellAnno, R. and Schneider, F. (2006). Estimating the Underground Economy by Using MIMIC Models: A Response to T. Breusch´s Critique. Technical report.

Giles, D. (1999). Measuring the Hidden Economy: Implications for Econometric Modelling. *The Economic Journal*, 109(456):370–380.

Glymour, C., Scheines, R., Spirtes, P., and Kelly, K. (1987). *Discovering Causal Structure: Artificial Intelligence, Philosophy of Science, and Statistical Modeling*. Academic Press.

Hair, J., Tatham, R., R.E., A., and Black, W. (1998). *Multivariate Data Analysis*. Prentice-Hall International.

Hastie, T., Tibshirani, R., Friedman, J., Hastie, T., Friedman, J., and Tibshirani, R. (2009). *The elements of statistical learning*, volume 2. Springer.

Hempel, C. G. (1985). Thoughts on the limitations of discovery by computer. *Logic of discovery and diagnosis in medicine*, pages 115–122.

Kant, I., Guyer, P., and Wood, A. W. (1998). *Critique of pure reason*. Cambridge University Press.

Kummerfeld, E., Ramsey, J., Yang, R., Spirtes, P., and Scheines, R. (2014). Causal clustering for 2-factor measurement models. *Proceedings of the European Conference on Machine Learning. (Forthcoming)*.

Lester, L. (2008). A Multiple Indicators and Multiple Causes (MIMIC) Model of Immigrant Settlement Success.

Logothetis, N. K. (2008). What we can do and what we cannot do with fmri. *Nature*, 453(7197):869–878.

Markowetz, F., Kostka, D., Troyanskaya, O. G., and Spang, R. (2007). Nested effects models for high-dimensional phenotyping screens. *Bioinformatics*, 23(13):i305–i312.

Mayo, D. G. (1996). *Error and the growth of experimental knowledge*. University of Chicago Press.

Network, C. G. A. R. et al. (2011). Integrated genomic analyses of ovarian carcinoma. *Nature*, 474(7353):609–615.

Ríos-Bedoya, C., Pomerleau, C., Neuman, R., and Pomerleau, O. (2009). Using MIMIC Models to Examine the Relationship Between Current Smoking and Early Smoking Experiences. *Nicotine & Tobacco Research*, 11(9):1035–1041.

Shalizi, C. R. (2012). Advanced data analysis from an elementary point of view. *Preprint of book found at http://www. stat. cmu. edu/ cshalizi/ADAfaEPoV*.

Silva, R., Scheine, R., Glymour, C., and Spirtes, P. (2006). Learning the Structure of Linear Latent Variable Models. *The Journal of Machine Learning Research*, 7:191–246.

Sober, E. (1998). Black Box Inference: When Should Intervening Variables be Postulated? *The British Journal for the Philosophy of Science*, 49(3):469–498.

Spirtes, P., Glymour, C., and Scheines, R. (2000). *Causation Prediction & Search 2e*. Adaptive Computation and Machine Learning Series. MIT Press.

Tedds, L. (1998). Measuring the Size of the Hidden Economy in Canada: A Latent Variable/MIMIC Model Approach. *MA Extended Essay, Department of Economics, University of Victoria*.

Thurstone, L. (1934). The Vectors of Mind. *Psychological Review; Psychological Review*, 41(1):1.

Tresch, A. and Markowetz, F. (2008). Structure learning in nested effects models. *Statistical Applications in Genetics and Molecular Biology*, 7(1).