

CAUSAL NETWORKS

CAUSAL DISCOVERY FROM OBSERVATIONAL DATA

Fabio Stella

Department of Informatics, Systems and Communication,

University of Milan-Bicocca

Viale Sarca 336, 2016 Milan, ITALY

e-mail: fabio.stella@unimib.it

Twitter: [FaSt@FabioAStella](https://twitter.com/FaSt@FabioAStella)

Previous lectures assumed that the causal graph is given. What if we don't know the graph? Can we learn it? We will refer to this problem as structure identification or structure learning. In this lecture we show how the structure of the causal network can be learnt from observational data.

In particular, the lecture presents the following:

- Constraint-based algorithms
- The PC algorithm
- Semi-parametric causal discovery
- Additional topics

PART I

CONSTRAINT-BASED METHODS

THE FUNDAMENTAL PROBLEM OF CAUSAL INFERENCE

We cannot observe both $Y_i(1)$ and $Y_i(0)$, therefore we cannot observe the causal effect

$$\tau_i \triangleq Y_i(1) - Y_i(0)$$

Another relevant task is that of **CAUSAL DISCOVERY TASK**, i.e., the problem of discovering the causal model which helps to **EXPLAIN** the causal effect of the treatment X on the outcome Y .

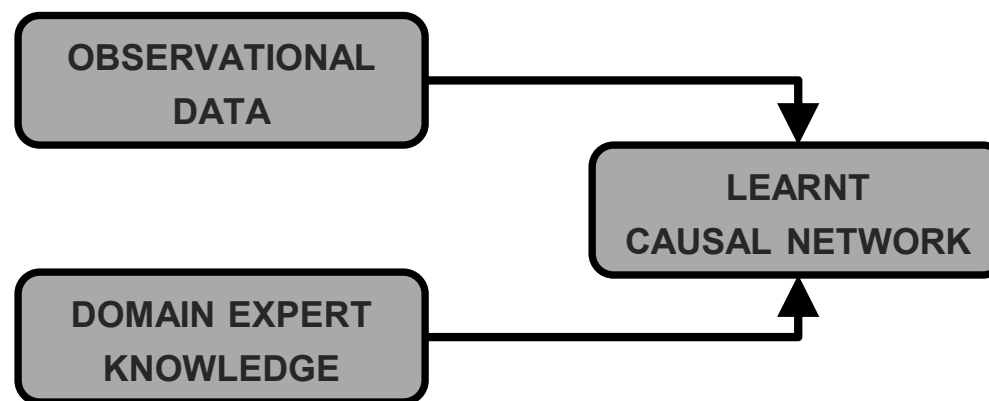


Figure 12.1

The **CAUSAL INFERENCE TASK** is an extremely relevant task and it consists of computing the causal effect τ_i of the treatment X on the outcome Y , no matter which is the causal model responsible for such a causal effect.

In this lecture we focus the attention to **CAUSAL NETWORKS**, as a valid tool to accomplish the **CAUSAL DISCOVERY TASK**.

More precisely we present the problem of **LEARNING A CAUSAL NETWORK** by fusion of **OBSERVATIONAL DATA** and **DOMAIN EXPERT KNOWLEDGE**.

Learning a **CAUSAL NETWORK** consists of learning

- Structure,
- Parameters.

Learning a **CAUSAL NETWORK** can be achieved by

- Constraint-based algorithms,
- Score-based algorithms,
- Hybrid algorithms.

However, in this lecture we focus the attention to **CONSTRAINT-BASED ALGORITHMS**.

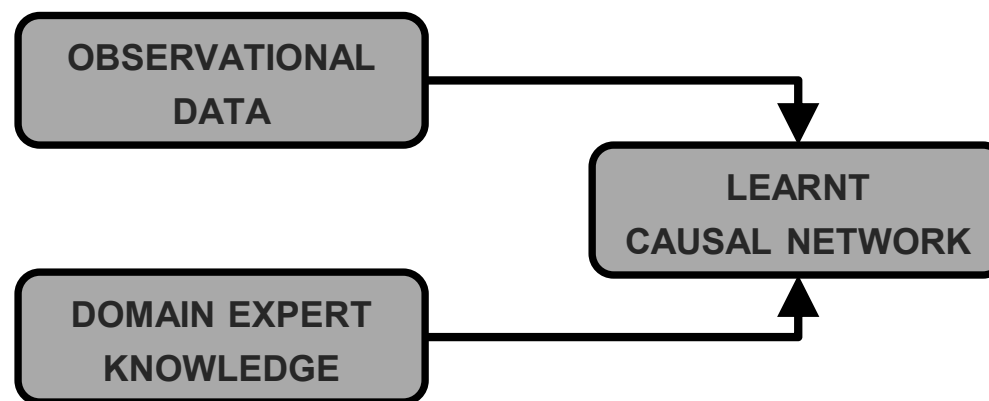


Figure 12.1

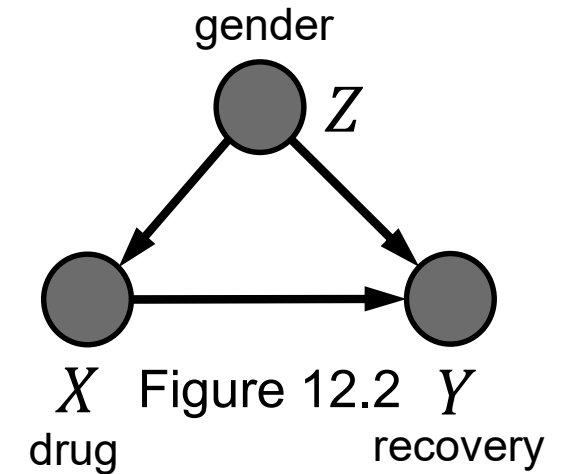


Figure 12.2

$$P(Z = 0) = 0.6$$

$$P(X = 0|Z = 0) = 0.5$$

$$P(X = 0|Z = 1) = 0.3$$

$$P(Y = 0|Z = 0, X = 0) = 0.6$$

$$P(Y = 0|Z = 1, X = 0) = 0.5$$

$$P(Y = 0|Z = 0, X = 1) = 0.1$$

$$P(Y = 0|Z = 1, X = 1) = 0.2$$

- We assume the **UNDERLYING PROCESS** follows a probability distribution P (the underlying probability distribution associated with DAG \mathcal{G}).
- Then, the **UNDERLYING PROCESS** can be adequately represented by sampling from P to obtain **OBSERVATIONAL DATA**.

The goal of the **CAUSAL DISCOVERY TASK** is to identify a model representation M of P .

To simplify the task, we assume the **PROBABILITY DISTRIBUTION P** to be a **DAG-FAITHFUL PROBABILITY DISTRIBUTION** with underlying DAG \mathcal{G} .

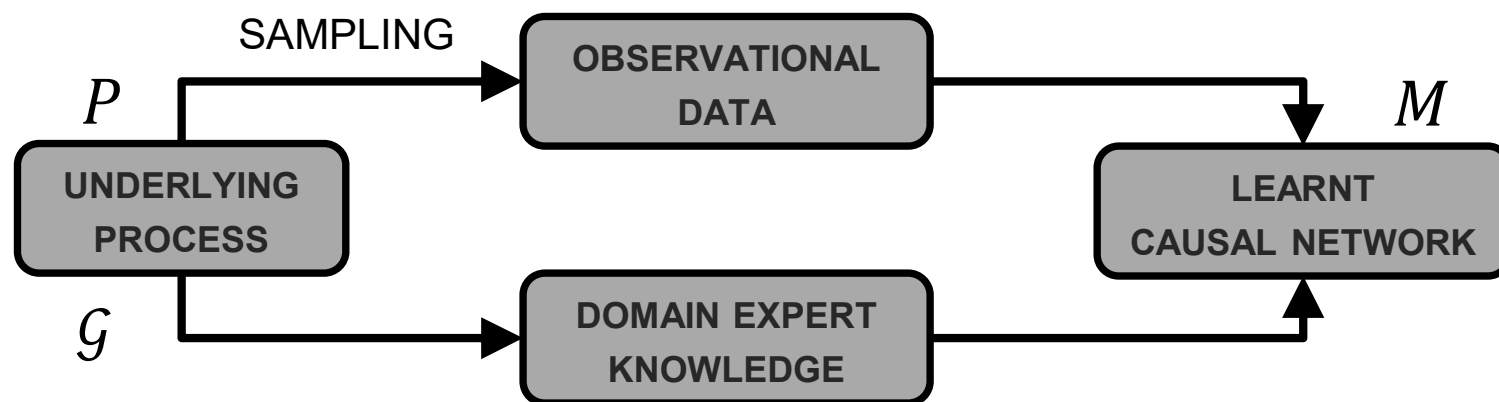


Figure 12.1

STABILITY – FAITHFULNESS

P is a stable (faithful) distribution if there exists a DAG \mathcal{G} such that

$$\mathbf{X} \perp\!\!\!\perp_P \mathbf{Y} \mid \mathbf{Z} \Leftrightarrow \mathbf{X} \perp\!\!\!\perp_{\mathcal{G}} \mathbf{Y} \mid \mathbf{Z}$$

for any three sets of variables \mathbf{X} , \mathbf{Y} , and \mathbf{Z} .

We assume that the distribution P can be represented as a **CAUSAL NETWORK** (if P is not **DAG-FAITHFUL**, a causal network may still be an excellent approximation).

The **FAITHFULNESS** assumption says that the distribution P , induced by \mathcal{G} , satisfies no independence relations beyond those implied by \mathcal{G} .

A **CAUSAL NETWORK** is **FAITHFUL** iff for every d-connection (no d-separation) \mathbf{Z} there is a corresponding conditional dependence, i.e.,

$$\mathbf{X} \not\perp\!\!\!\perp_{\mathcal{G}} \mathbf{Y} \mid \mathbf{Z} \Rightarrow \mathbf{X} \not\perp\!\!\!\perp_P \mathbf{Y} \mid \mathbf{Z}$$

FAITHFULNESS is a much less attractive assumption than the **MARKOV ASSUMPTION** because it is easy to think of counterexamples where

- two variables are independent in P ,
- but there are unblocked paths between them in \mathcal{G} .

It is worthwhile to mention that many **CONSTRAINT-BASED METHODS** also assume that there are no unobserved confounders, which is known as **CAUSAL SUFFICIENCY**.

CAUSAL SUFFICIENCY

There are no unobserved confounders of any of the variables in the graph.

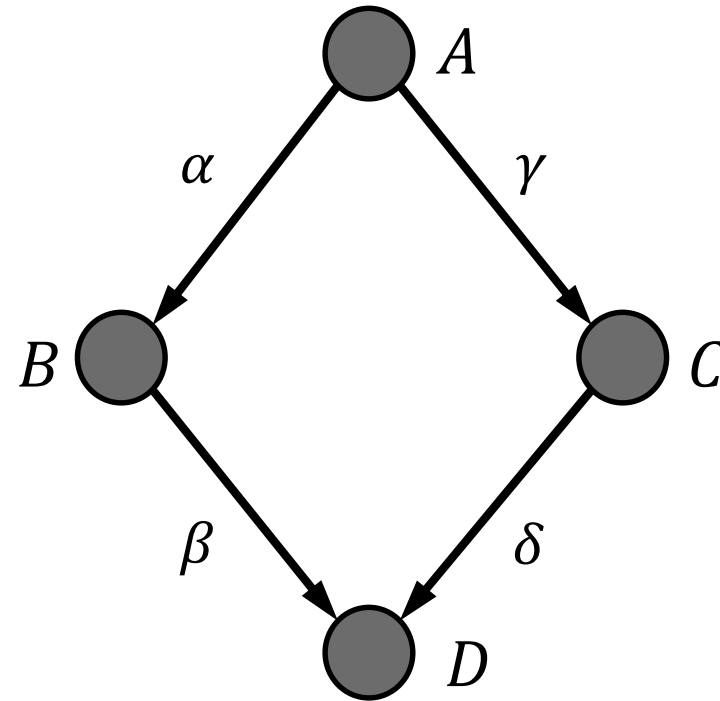


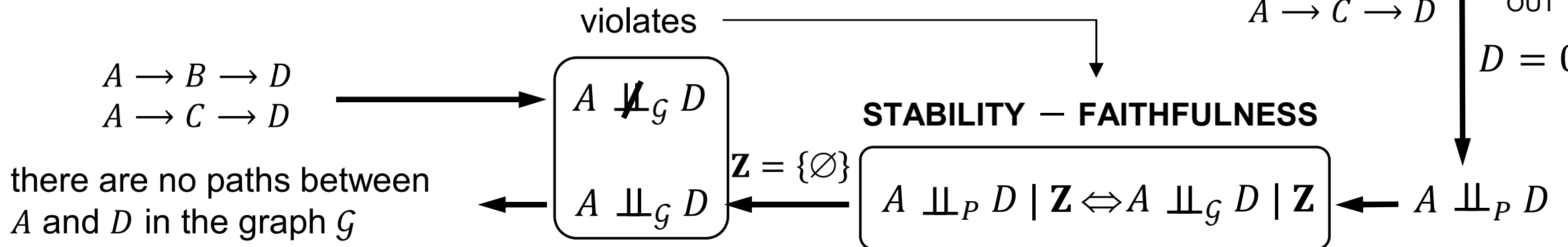
Figure 12.3

$$\begin{aligned}
 B &:= \alpha A \\
 C &:= \gamma A \\
 D &:= \beta B + \delta C
 \end{aligned}$$

$$\begin{aligned}
 D &= \beta\alpha A + \delta\gamma A \\
 D &= (\beta\alpha + \delta\gamma)A
 \end{aligned}$$

If $\beta\alpha = -\delta\gamma$

$$\begin{array}{l}
 A \rightarrow B \rightarrow D \\
 A \rightarrow C \rightarrow D
 \end{array}
 \begin{array}{l}
 \text{CANCEL} \\
 \text{OUT} \\
 D = 0
 \end{array}$$



cases generated by the underlying and unknown process
 $\mathcal{D} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$
 independent and identically distributed data cases drawn at random from the probability distribution P

assignment of values to the n variables for the j^{th} data case

$$\mathbf{x}^j = \{x_1^j, x_2^j, \dots, x_n^j\}$$

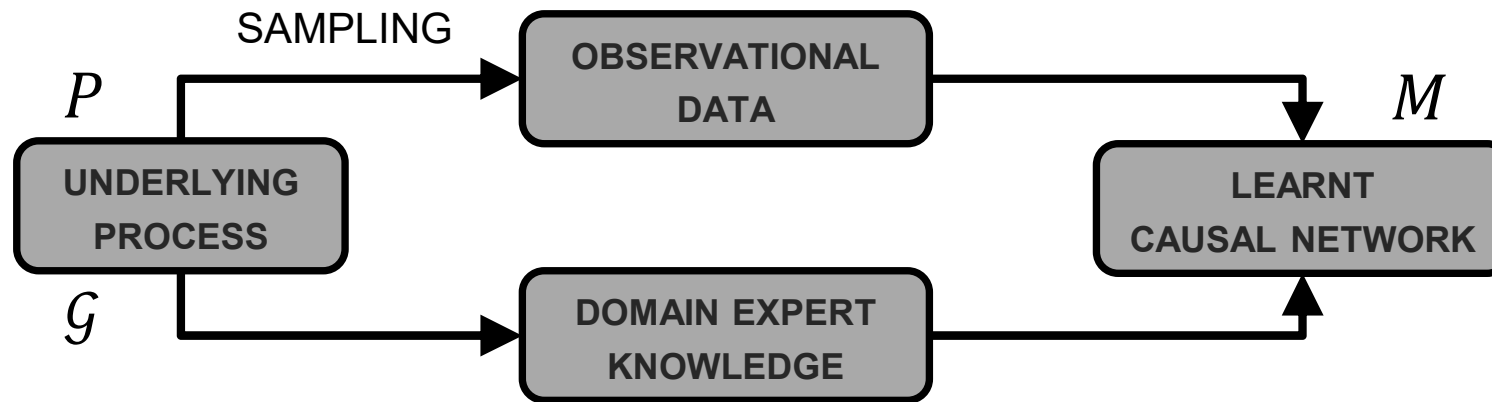


Figure 12.1

Some values in \mathcal{D} may be missing (N/A), but missing values are assumed to be:

- **MISSING AT RANDOM (MAR)** or,
- **MISSING COMPLETELY AT RANDOM (MCAR)**.
 i.e., the missing data mechanism is uninformative and can be ignored.

A variable never observed is called a **HIDDEN** or a **LATENT VARIABLE**.

LEARNING A CAUSAL NETWORK is the **TASK** of identifying a DAG structure \mathcal{G} and a set of conditional probability distributions with parameters Θ on the basis of

- $\mathcal{D} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$ and,
- possibly some domain expert background knowledge.

	X_1	X_2	...	X_n
\mathbf{x}^1	blue	yes	...	low
\mathbf{x}^2	green	no	...	low
\mathbf{x}^3	red	N/A	...	high
...

Table 12.1

Article

Hard and Soft EM in Bayesian Network Learning from Incomplete Data

Andrea Ruggieri ^{1,†}, Francesco Stranieri ^{1,†}, Fabio Stella ¹ and Marco Scutari ^{2,*}

¹ Department of Informatics, Systems and Communication, Università degli Studi di Milano-Bicocca, 20126 Milano MI, Italy; a.ruggieri4@campus.unimib.it (A.R.); f.stranieri1@campus.unimib.it (F.S.); fabio.stella@unimib.it (F.S.)

² Istituto Dalle Molle di Studi sull'Intelligenza Artificiale (IDSIA), 6962 Viganello, Switzerland

* Correspondence: scutari@idsia.ch

† These authors contributed equally to this work.

Received: 18 November 2020; Accepted: 7 December 2020; Published: 9 December 2020



Abstract: Incomplete data are a common feature in many domains, from applications. Bayesian networks (BNs) are often used in these domains because of their graphical and causal interpretations. BN parameter learning from incomplete data is usually implemented with the Expectation-Maximisation algorithm (EM), which computes the relevant sufficient statistics (“soft EM”) using belief propagation. Similarly, the Structural Expectation-Maximisation algorithm (Structural EM) learns the network structure of the BN from those sufficient statistics using algorithms designed for complete data. However, practical implementations of parameter and structure learning

In the **CONSTRAINT-BASED APPROACH**, the DAG \mathcal{G} of a causal network is considered as an encoding of a set of (conditional) dependence and independence relations $\mathcal{M}_{\mathcal{G}}$, which can be read off \mathcal{G} using **D-SEPARATION**.

D-SEPARATION

A path p is blocked by a set of nodes S if and only if

- 1) p contains a chain of nodes $A \rightarrow B \rightarrow C$ or a fork $A \leftarrow B \rightarrow C$ such that the middle node B is in S (i.e., is conditioned on),
- 2) or p contains a collider $A \rightarrow B \leftarrow C$ such that the collision node B is not in S , and no descendant of B is in S .

If S blocks every path between two nodes X and Y , then X and Y are d-separated, conditional on S , and thus are independent conditional on S .

- Structure learning is then the task of identifying a DAG structure that (best) encodes a set of (conditional) dependence and independence relations $\mathcal{M}_{\mathcal{G}}$.
- The set of (conditional) dependence and independence relations $\mathcal{M}_{\mathcal{G}}$ may, for instance, be derived from **OBSERVATIONAL DATA** by statistical tests.
- Based on \mathfrak{D} alone, we can at most hope to identify an **EQUIVALENCE CLASS OF GRAPHS** encoding the (conditional) dependence and independence relations $\mathcal{M}_{\mathcal{G}}$ of the generating distribution P .

A **CONSTRAINT-BASED STRUCTURE LEARNING ALGORITHM** proceeds by determining the validity of independence relations of the form:

$I(X, Y | S_{XY})$ we check whether X is independent of Y given subset S_{XY} , where $X, Y \in \mathbf{X}$ and $S_{XY} \subseteq \mathbf{X}$

- The structure learning algorithm will work with any information source able to provide such information.
- We will consider the case where the validity of independence relations is determined by **STATISTICAL HYPOTHESIS TESTS OF INDEPENDENCE** based on a database of cases (**OBSERVATIONAL DATA**).

The size of the space of possible DAGs grows **SUPER-EXPONENTIALLY** with the number of nodes n in the graph.

The following recursive formula gives the number $f(n)$ of DAGs on n nodes:

$$f(n) = \sum_{i=1}^n (-1)^{i+1} \frac{n!}{(n-i)! i!} 2^{i(n-1)} f(n-i)$$

n	#dags
1	1
2	3
3	25
4	543
5	29,281
6	3,781,503
7	$1.1 \cdot 10^9$
8	$7.8 \cdot 10^{11}$
9	$1.2 \cdot 10^{15}$
10	$4.2 \cdot 10^{18}$

Table 12.2

Under the conditions listed below:

- 1) The independence relationships have a perfect representation as a DAG (**ACYCLICITY AND FAITHFULNESS ASSUMPTIONS**).
- 2) No hidden (latent) variables are involved (**CAUSAL SUFFICIENCY ASSUMPTION**).
- 3) The database (**OBSERVATIONAL DATA**) consists of a set of independent and identically distributed cases.
- 4) The database (**OBSERVATIONAL DATA**) is infinitely large.
- 5) The statistical tests have no error.

a constraint-based structure learning algorithm discovers a DAG structure equivalent to the DAG structure of P .

- Two DAGs representing the same set of (conditional) dependence and independence relations are equivalent in the sense that they can capture the same set of probability distributions.
- The equivalence class of a DAG \mathcal{G} is the set of DAGs with the same set of d-separation relations as \mathcal{G} .
- A PDAG — an acyclic, partially directed graph, i.e., an acyclic graph with some edges undirected (also known as a pattern or **ESSENTIAL GRAPH**) — can be used to represent the equivalence class of a set of DAG structures, i.e., a maximal set of DAGs with the same set of d-separation relations.

EQUIVALENT MODELS AND MARKOV EQUIVALENCE CLASS

Any two models $M_{\mathcal{G}'}$ and $M_{\mathcal{G}''}$ over the same set of variables, whose graphs \mathcal{G}' and \mathcal{G}'' , respectively, have the same skeleton \mathcal{G}_S (i.e., undirected graph obtained by replacing directed edges with undirected edges) and the same v-structures, are **EQUIVALENT**.

Two graphs \mathcal{G}' and \mathcal{G}'' are in the same **EQUIVALENCE CLASS**

- if they share a **COMMON SKELETON**—that is, if they possess the same edges, regardless of the direction of those edges—and
- if they share **COMMON V-STRUCTURES**, that is, colliders whose parents are not adjacent.

The three graphs in Figure 12.4 are **MARKOV EQUIVALENT**.

Hence, based on data $\mathcal{D} = \{x^1, x^2, \dots, x^N\}$ alone we cannot distinguish between these three models, while we can distinguish them from the collider in Figure 12.6.

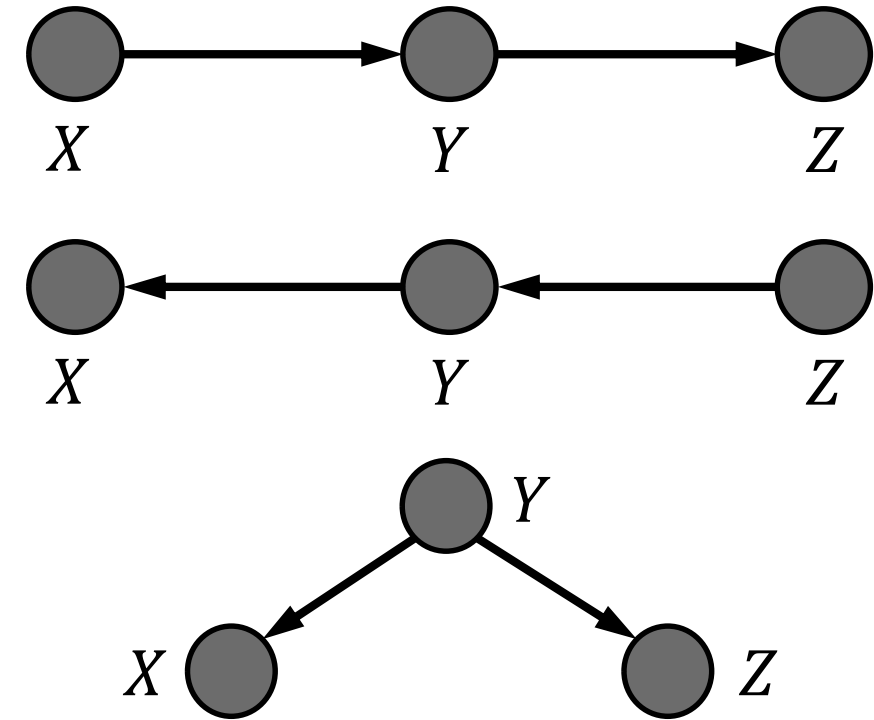


Figure 12.4

- common skeleton
- no v-structures

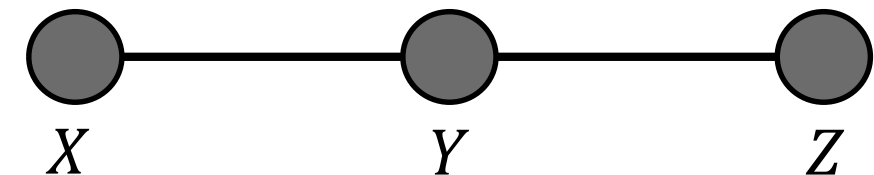


Figure 12.5

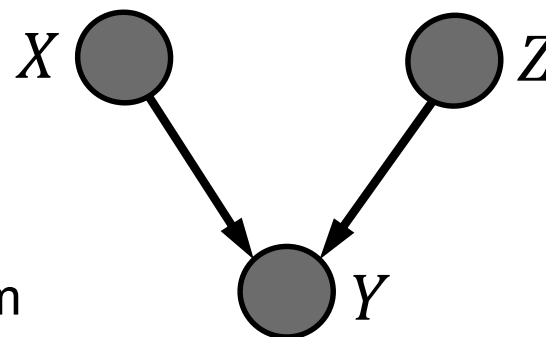


Figure 12.6

The graphs in Figure 12.4 correspond to the same set of independence/dependence assumptions \mathcal{M}_G , i.e., we say that the three graphs are **MARKOV EQUIVALENT**.

Given a graph, we refer to its **MARKOV EQUIVALENCE CLASS** as the set of graphs that encode the same (conditional) independencies.

- $X \perp\!\!\!\perp Z$ X is independent on Z
- $X \not\perp\!\!\!\perp Z$ X is not independent on Z
- $X \perp\!\!\!\perp Z \mid Y$ X is independent on Z given Y
- $X \not\perp\!\!\!\perp Z \mid Y$ X is not independent on Z given Y

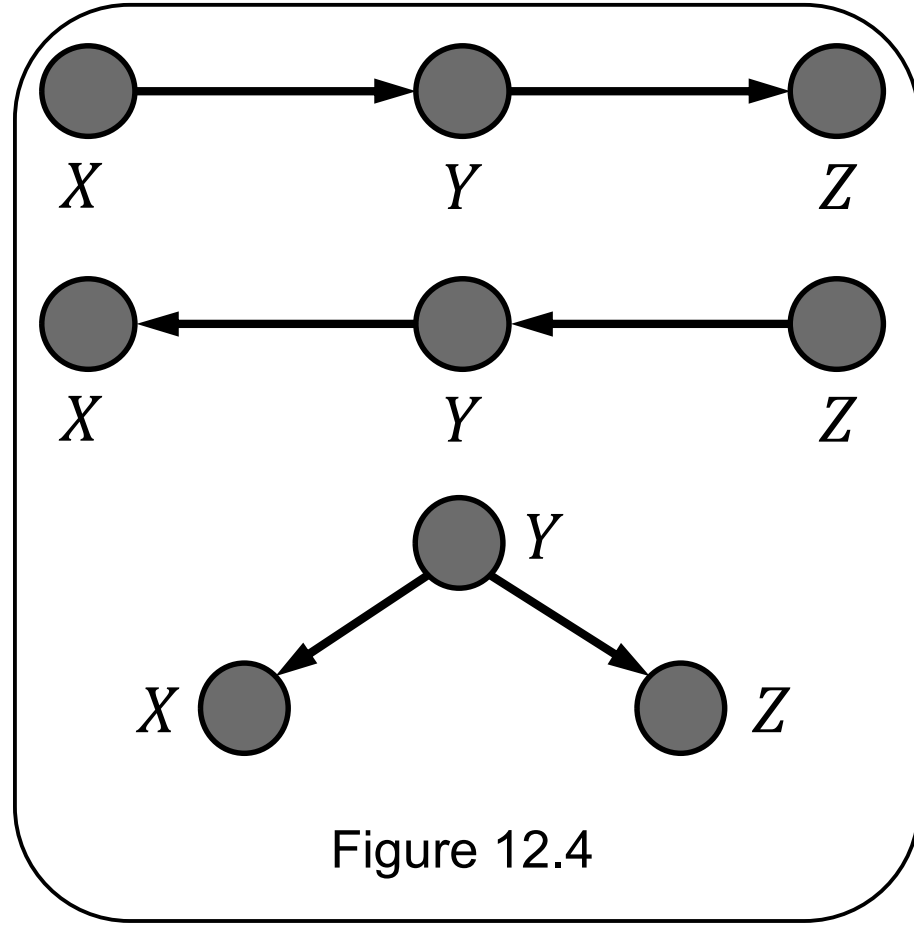
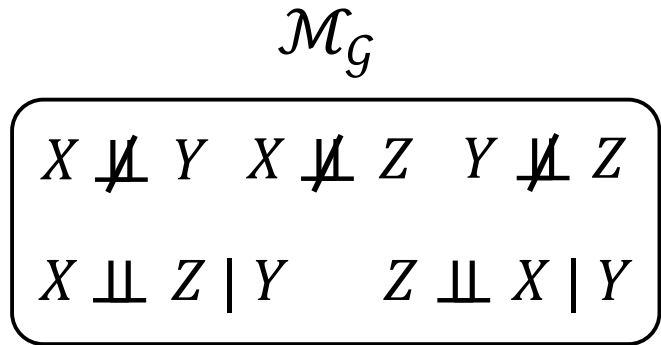
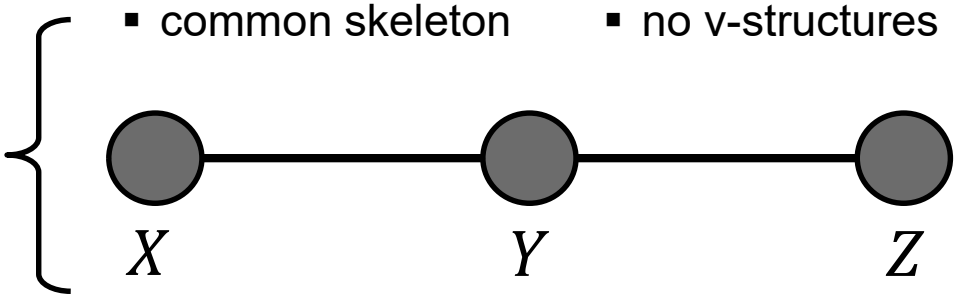


Figure 12.4

The three graphs can not be distinguished, when all we have is observational data!!!

The three graphs are **MARKOV EQUIVALENT**



EQUIVALENCE CLASS

An equivalence class is a maximal set of DAGs with the same set of independence properties \mathcal{M}_G .

The three DAGs in Figure 12.7 all represent the same set of conditional independence and dependence relations.

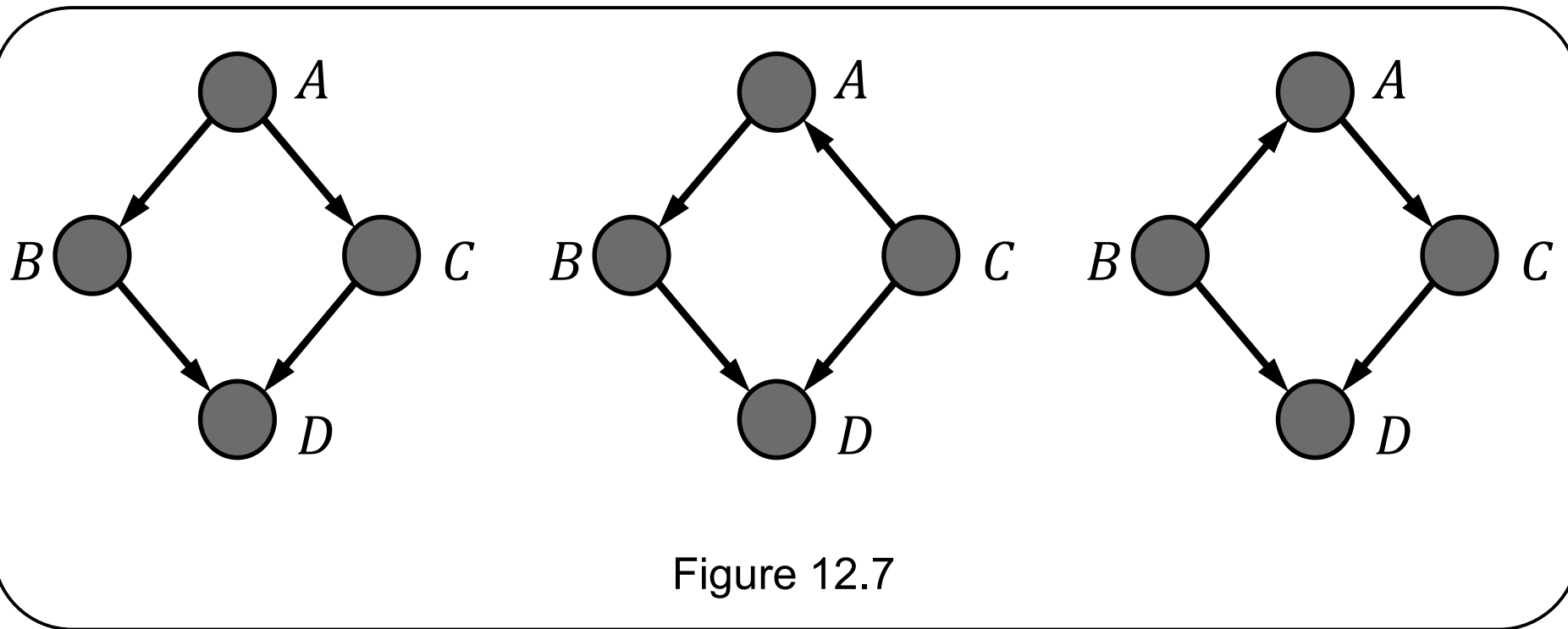


Figure 12.7

Three equivalent DAGs

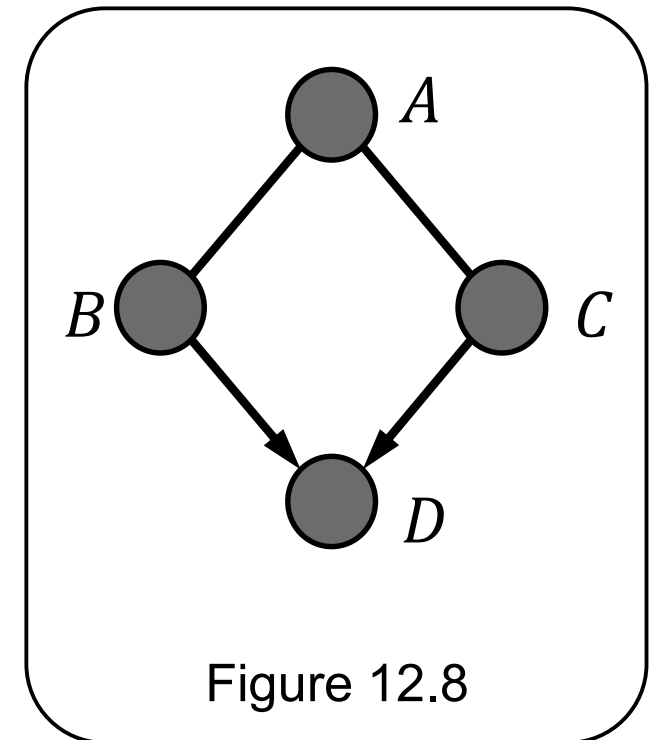


Figure 12.8

Equivalence class of the three DAGs to the left

A set of (conditional) dependence and independence relations \mathcal{M}_G may be generated by statistical tests on the **OBSERVATIONAL DATA**.

In each test, the hypothesis tested is that of independence between a pair of variables.

Let X and Y be a pair of variables for which we would like to determine dependence by **STATISTICAL HYPOTHESIS TESTING**.

We could:

- test for **MARGINAL INDEPENDENCE** and subsequently,
- test for **CONDITIONAL INDEPENDENCE** given subsets of other variables.

	X_1	X_2	...	X_n
x^1	blue	yes	...	low
x^2	green	no	...	low
x^3	red	N/A	...	high
...

Table 12.1

A set of (conditional) dependence and independence relations \mathcal{M}_G may be generated by statistical tests on the **OBSERVATIONAL DATA**.

In each test, the hypothesis tested is that of independence between a pair of variables.

Let X and Y be a pair of variables for which we would like to determine dependence by **STATISTICAL HYPOTHESIS TESTING**.

We could:

- test for **MARGINAL INDEPENDENCE** and subsequently,
- test for **CONDITIONAL INDEPENDENCE** given subsets of other variables.

In the case of **MARGINAL INDEPENDENCE TESTING** between X and Y , the **HYPOTHESIS TO BE TESTED** is

$$H_0 : P(X, Y) = P(X) P(Y) \quad X \perp\!\!\!\perp_P Y$$

$$H_1 : P(X, Y) \neq P(X) P(Y)$$

	X_1	X_2	...	X_n
x^1	blue	yes	...	low
x^2	green	no	...	low
x^3	red	N/A	...	high
...

Table 12.1

A potential hypothesis test

Under H_0 the likelihood statistic

$$G^2 = 2 \sum_{x,y} N_{xy} \log \left(\frac{N_{xy}}{\mathbb{E}_{xy}} \right) \quad \mathbb{E}_{xy} = \frac{N_x N_y}{N}$$

has an asymptotic χ^2 distribution with the appropriate degrees of freedom (df).

N_{xy} , number of cases where $X = x$ and $Y = y$

N_x , number of cases where $X = x$

N_y , number of cases where $Y = y$

A set of (conditional) dependence and independence relations \mathcal{M}_G may be generated by statistical tests on the **OBSERVATIONAL DATA**.

In each test, the hypothesis tested is that of independence between a pair of variables.

Let X and Y be a pair of variables for which we would like to determine dependence by **STATISTICAL HYPOTHESIS TESTING**.

We could:

- test for **MARGINAL INDEPENDENCE** and subsequently,
- test for **CONDITIONAL INDEPENDENCE** given subsets of other variables.

In the case of **CONDITIONAL INDEPENDENCE TESTING** between X and Y given subset S_{XY} the **HYPOTHESIS TO BE TESTED** is

$$H_0 : P(X, Y | S_{XY}) = P(X | S_{XY}) P(Y | S_{XY}) \quad X \perp\!\!\!\perp Y | S_{XY}$$

$$H_1 : P(X, Y | S_{XY}) \neq P(X | S_{XY}) P(Y | S_{XY})$$

	X_1	X_2	...	X_n
x^1	blue	yes	...	low
x^2	green	no	...	low
x^3	red	N/A	...	high
...

Table 12.1

A potential hypothesis test

Under H_0 the likelihood statistic

$$G^2 = 2 \sum_{x,y,s} N_{xys} \log \left(\frac{N_{xys}}{\mathbb{E}_{xys}} \right) \quad \mathbb{E}_{xys} = \frac{N_{xs} N_{ys}}{N_s}$$

has an asymptotic χ^2 distribution with the appropriate degrees of freedom (df).

$$df = (\|X\| - 1)(\|Y\| - 1) \prod_{s \in S_{XY}} \|Y\|$$

PART II

THE PC ALGORITHM

The **PC ALGORITHM** (Spirtes & Glymour 1991, Spirtes et al. 2000) is probably the most known **CONSTRAINT-BASED ALGORITHM** for learning the structure of a causal network.

The **MAIN STEPS OF THE PC ALGORITHM** are:

- 1) Test for (conditional) independence between each pair of variables represented in $\mathcal{D} = \{x^1, x^2, \dots, x^N\}$ to derive $\mathcal{M}_{\mathcal{D}}$, the set of conditional independence and dependence relations.
- 2) Identify the skeleton of the graph induced by $\mathcal{M}_{\mathcal{D}}$.
- 3) Identify colliders.
- 4) Identify derived directions.

The **PC ALGORITHM** typically produces a **PDAG** (Partially DAG) representing an equivalence class as it emerges from hypothesis testing performed by using the available observational data $\mathcal{D} = \{x^1, x^2, \dots, x^N\}$.

STEP 1: TEST FOR (CONDITIONAL) INDEPENDENCE. We try to determine the validity of the conditional independence statement

$$X \perp\!\!\!\perp_P Y \mid S_{XY} \quad \begin{aligned} H_0 &: P(X, Y \mid S_{XY}) = P(X \mid S_{XY})P(Y \mid S_{XY}) \\ H_1 &: P(X, Y \mid S_{XY}) \neq P(X \mid S_{XY})P(Y \mid S_{XY}) \end{aligned}$$

The independence hypothesis H_0 is tested for **CONDITIONING SETS** S_{XY} of cardinality 0, 1, 2, 3, ... in that order.

If the hypothesis H_0 cannot be rejected based on some preselected **SIGNIFICANCE LEVEL** α , then the search for an independence relation between X and Y is terminated.

Assume you are given the database of cases (**OBSERVATIONAL DATA**) which has been generated from the model in Figure 12.9, and that you got the following (conditional) independence and dependence relations:

$$\mathcal{M}_{\mathcal{D}} \left\{ \begin{aligned} \mathcal{M}_{\perp} &= \{B \perp\!\!\!\perp E, B \perp\!\!\!\perp R, B \perp\!\!\!\perp W \mid A, A \perp\!\!\!\perp R \mid E, E \perp\!\!\!\perp W \mid A, R \perp\!\!\!\perp W \mid A\} \leftarrow \text{independencies} \\ \mathcal{M}_{\not\perp} &= \{B \not\perp\!\!\!\perp A, B \not\perp\!\!\!\perp A \mid \{E\}, B \not\perp\!\!\!\perp A \mid \{R\}, B \not\perp\!\!\!\perp A \mid \{W\}, B \not\perp\!\!\!\perp A \mid \{E, R\}, \\ &\quad B \not\perp\!\!\!\perp A \mid \{E, W\}, B \not\perp\!\!\!\perp A \mid \{R, W\}, B \not\perp\!\!\!\perp A \mid \{E, R, W\}, A \not\perp\!\!\!\perp E, \dots \\ &\quad A \not\perp\!\!\!\perp W, \dots, E \not\perp\!\!\!\perp R, \dots\}. \leftarrow \text{dependencies} \end{aligned} \right.$$

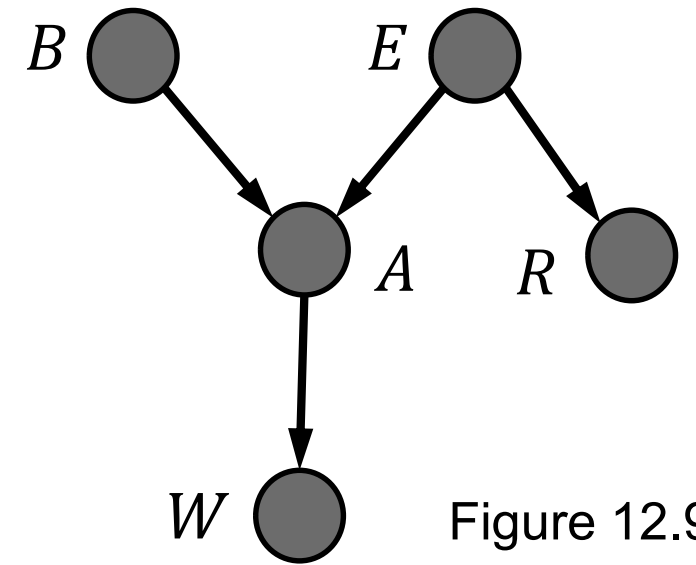


Figure 12.9

STEP 2: IDENTIFY THE SKELETON. The skeleton of an acyclic, directed or partially directed graph \mathcal{G} is the undirected graph \mathcal{G}^u obtained from \mathcal{G} by removing the direction on all directed edges.

The skeleton of the graph induced from $\mathcal{M}_{\mathcal{D}}$ is constructed from the conditional dependence and independence statements of $\mathcal{M}_{\mathcal{D}}$ generated by the statistical test in **STEP 1**.

For each pair of variables X and Y where no independence statement

$$X \perp\!\!\!\perp_P Y \mid S_{XY}$$

exists, the undirected edge (X, Y) is created in the skeleton.

The graph of Figure 12.10 is a more intuitive and compact representation of the dependence and independence model than that of equations below.

$$\mathcal{M}_{\mathcal{D}} \left\{ \begin{array}{l} \mathcal{M}_{\perp} = \{B \perp E, B \perp R, B \perp W \mid A, A \perp R \mid E, E \perp W \mid A, R \perp W \mid A\} \\ \mathcal{M}_{\not\perp} = \{B \not\perp A, B \not\perp A \mid \{E\}, B \not\perp A \mid \{R\}, B \not\perp A \mid \{W\}, B \not\perp A \mid \{E, R\}, \\ B \not\perp A \mid \{E, W\}, B \not\perp A \mid \{R, W\}, B \not\perp A \mid \{E, R, W\}, A \not\perp E, \dots \\ A \not\perp W, \dots, E \not\perp R, \dots\}. \end{array} \right.$$

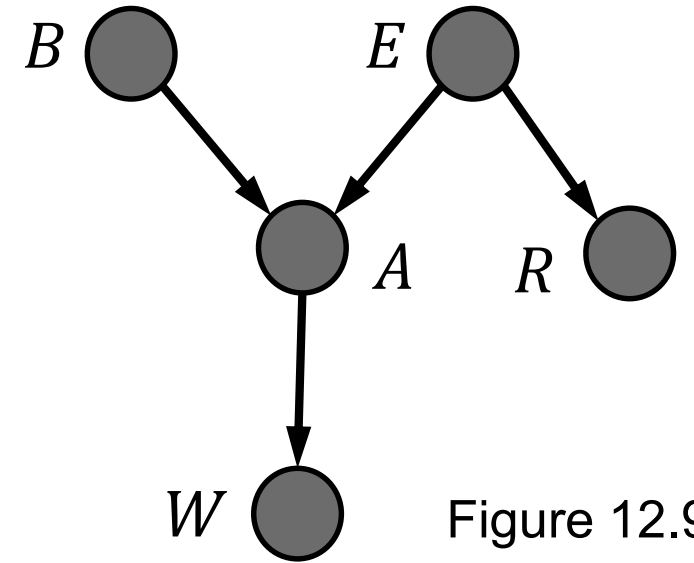


Figure 12.9

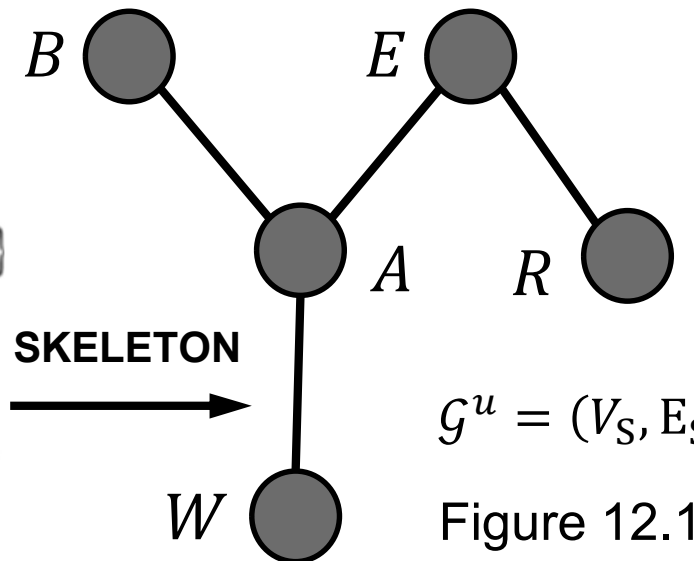


Figure 12.10

STEP 3: IDENTIFY COLLIDERS. Once the skeleton has been identified, colliders in the skeleton are identified.

Based on the skeleton, we search for subsets of variables $\{X, Y, Z\}$ such that X and Y are neighbors, Z and Y are neighbors while X and Z are not neighbors.

For each such subset a collider $X \rightarrow Y \leftarrow Z$ is created when $Y \notin S_{XZ}$ for any S_{XZ} satisfying

$$B \rightarrow A \leftarrow E$$

$$X \perp\!\!\!\perp_P Z \mid S_{XZ}$$

in $\mathcal{M}_{\mathcal{D}}$.

$$(B, A), (E, A) \in E_S$$

$$(B, E) \notin E_S$$

$$\boxed{\forall B \perp\!\!\!\perp_P E \mid S_{BE}} \quad S_{BE} = \{\emptyset\} \quad A \notin S_{BE}$$

$$\mathcal{M}_{\mathcal{D}} \left\{ \begin{array}{l} \mathcal{M}_{\perp} = \boxed{B \perp\!\!\!\perp E} B \perp\!\!\!\perp R, B \perp\!\!\!\perp W \mid A, A \perp\!\!\!\perp R \mid E, E \perp\!\!\!\perp W \mid A, R \perp\!\!\!\perp W \mid A \\ \mathcal{M}_{\not\perp} = \{B \not\perp\!\!\!\perp A, B \not\perp\!\!\!\perp A \mid \{E\}, B \not\perp\!\!\!\perp A \mid \{R\}, B \not\perp\!\!\!\perp A \mid \{W\}, B \not\perp\!\!\!\perp A \mid \{E, R\}, \\ B \not\perp\!\!\!\perp A \mid \{E, W\}, B \not\perp\!\!\!\perp A \mid \{R, W\}, B \not\perp\!\!\!\perp A \mid \{E, R, W\}, A \not\perp\!\!\!\perp E, \dots \\ A \not\perp\!\!\!\perp W, \dots, E \not\perp\!\!\!\perp R, \dots\}. \end{array} \right.$$

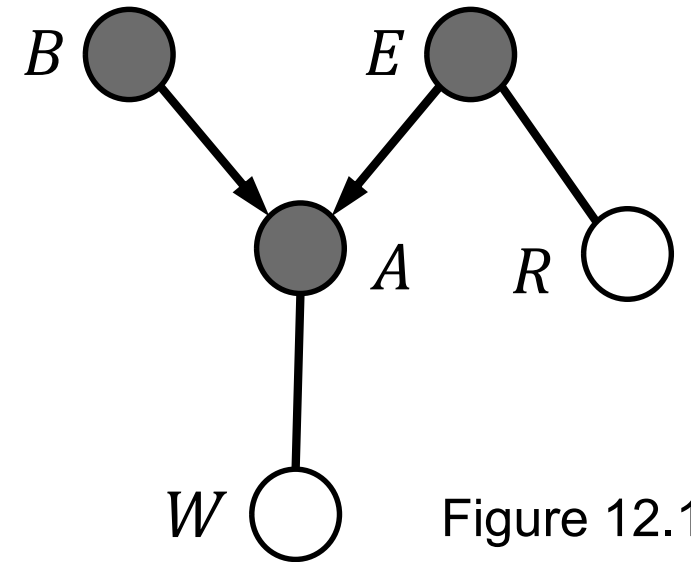


Figure 12.11

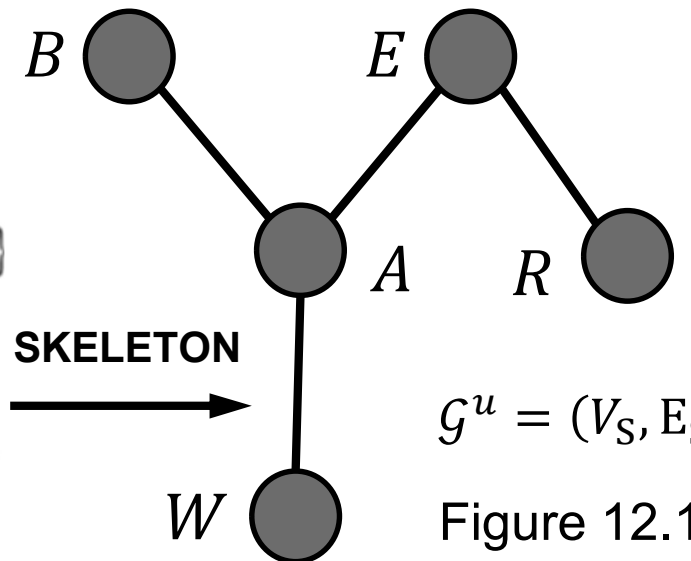


Figure 12.10

STEP 3: IDENTIFY COLLIDERS. Once the skeleton has been identified, colliders in the skeleton are identified.

Based on the skeleton, we search for subsets of variables $\{X, Y, Z\}$ such that X and Y are neighbors, Z and Y are neighbors while X and Z are not neighbors.

For each such subset a collider $X \rightarrow Y \leftarrow Z$ is created when $Y \notin S_{XZ}$ for any S_{XZ} satisfying

$$B \rightarrow A \leftarrow W$$

$$X \perp\!\!\!\perp_P Z \mid S_{XZ}$$

in $\mathcal{M}_{\mathcal{D}}$.

$$(B, A), (W, A) \in E_S$$

$$(B, W) \notin E_S$$

$$\boxed{\forall B \perp\!\!\!\perp_P W \mid S_{BW}} \quad S_{BW} = \{A\} \quad A \in S_{BW}$$

$$\mathcal{M}_{\mathcal{D}} \left\{ \begin{array}{l} \mathcal{M}_{\perp} = \{B \perp\!\!\!\perp E, B \perp\!\!\!\perp R, \boxed{B \perp\!\!\!\perp W \mid A}, A \perp\!\!\!\perp R \mid E, E \perp\!\!\!\perp W \mid A, R \perp\!\!\!\perp W \mid A\} \\ \mathcal{M}_{\not\perp} = \{B \not\perp\!\!\!\perp A, B \not\perp\!\!\!\perp A \mid \{E\}, B \not\perp\!\!\!\perp A \mid \{R\}, B \not\perp\!\!\!\perp A \mid \{W\}, B \not\perp\!\!\!\perp A \mid \{E, R\}, \\ B \not\perp\!\!\!\perp A \mid \{E, W\}, B \not\perp\!\!\!\perp A \mid \{R, W\}, B \not\perp\!\!\!\perp A \mid \{E, R, W\}, A \not\perp\!\!\!\perp E, \dots \\ A \not\perp\!\!\!\perp W, \dots, E \not\perp\!\!\!\perp R, \dots\}. \end{array} \right.$$

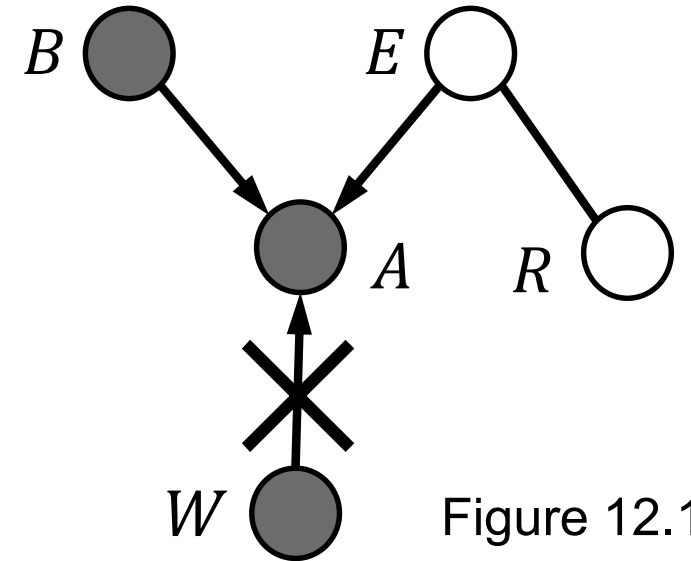
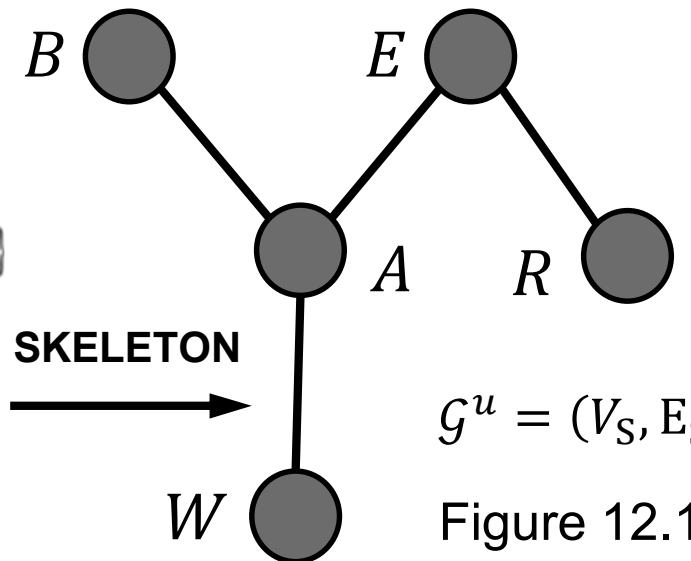


Figure 12.12



$$\mathcal{G}^u = (V_S, E_S)$$

Figure 12.10

STEP 4: IDENTIFY DERIVED DIRECTIONS. After identifying the skeleton and the colliders of \mathcal{G} , derived directions are identified.

The direction of an edge is said to be derived when it is a logical consequence of (the lack of) previous actions (i.e., since the edge was not directed in a previous step and it should have been in order to have a certain direction, then the edge must be directed in the opposite direction).

Starting with any PDAG including all valid colliders, a **MAXIMALLY DIRECTED PDAG** can be obtained following four **NECESSARY AND SUFFICIENT RULES**.

That is, by repeated application of these four rules all edges common to the equivalence class of \mathcal{G} are identified.

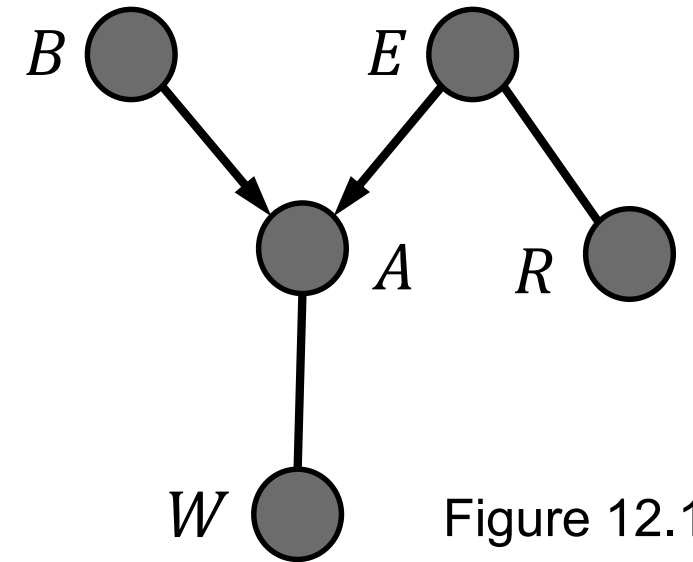
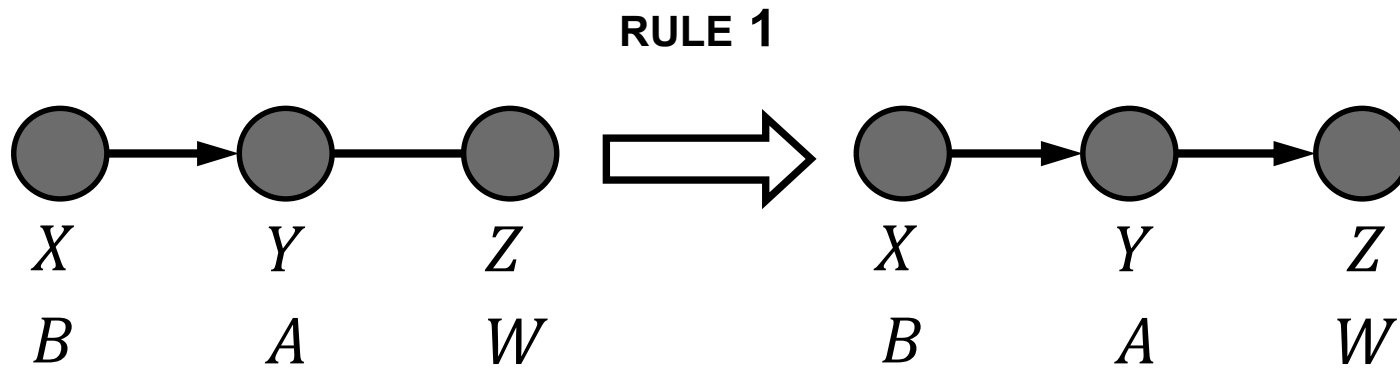


Figure 12.13

STEP 4: IDENTIFY DERIVED DIRECTIONS.



It follows from the fact that the collider

$$X \rightarrow Y \leftarrow Z$$

was not identified as a valid collider.

Since the edge between Y and Z is not part of the aforementioned collider, it must be directed from Y to Z .

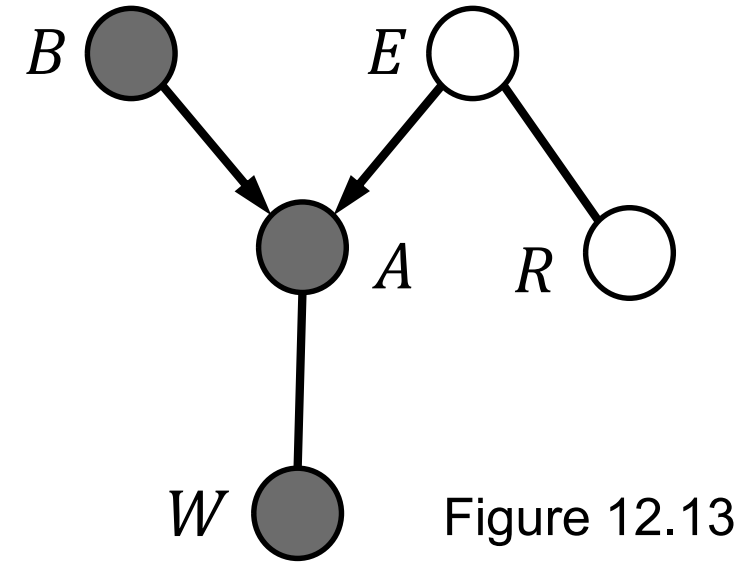


Figure 12.13

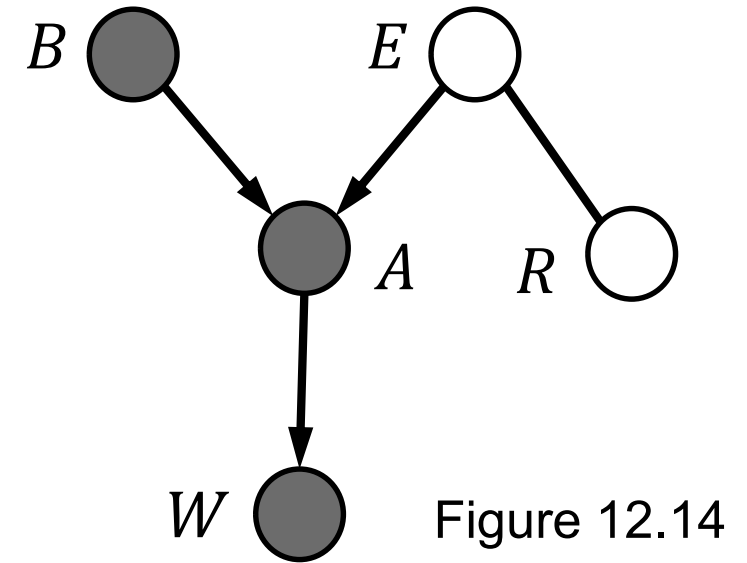
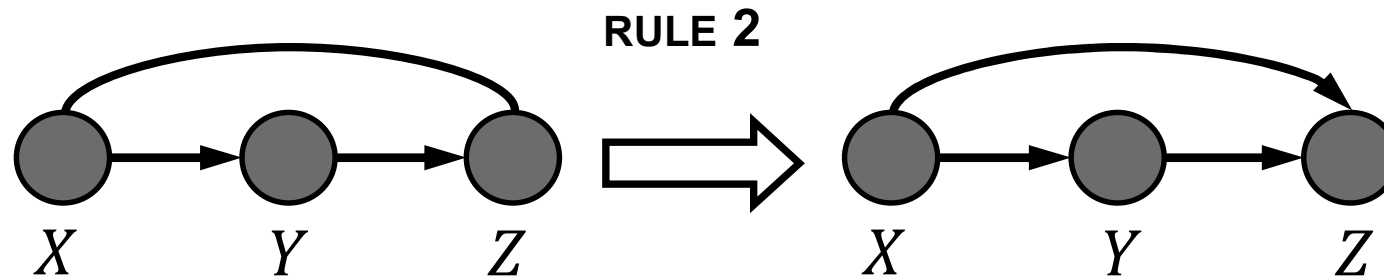


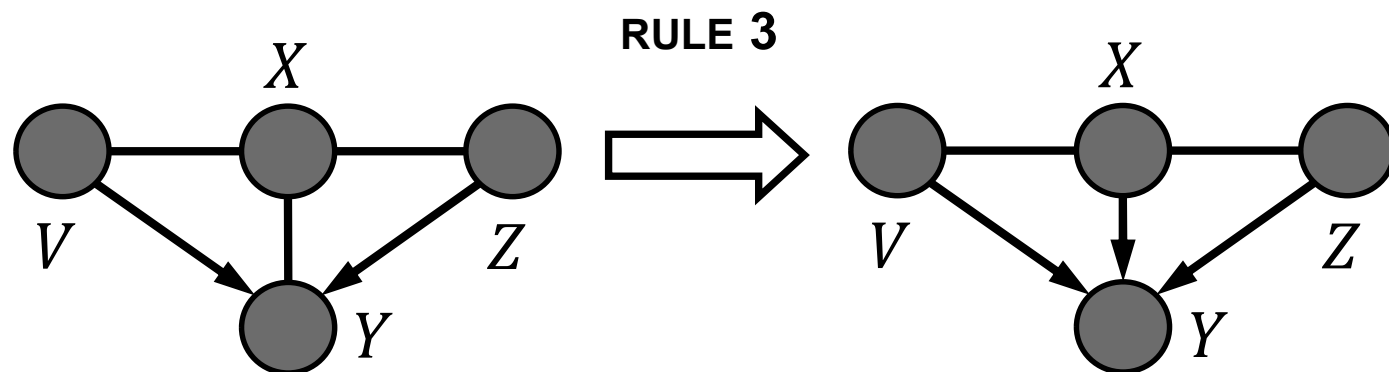
Figure 12.14

STEP 4: IDENTIFY DERIVED DIRECTIONS.



It follows from the fact that directing the edge between X and Z from Z to X will induce a directed cycle in the graph.

Thus, the edge must be directed from X to Z .

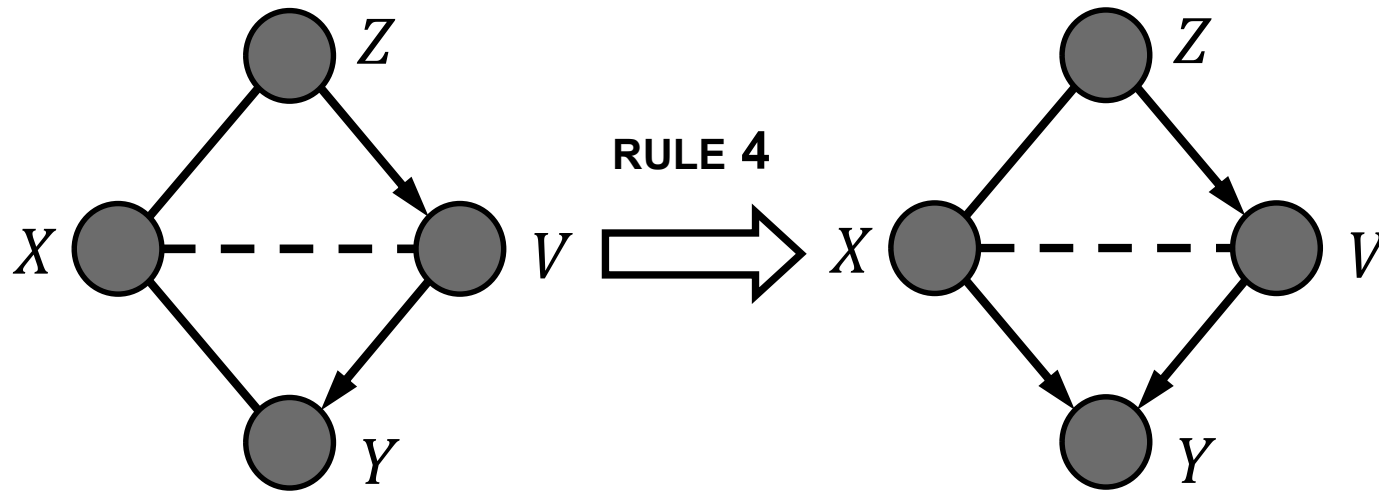


It follows from the fact that directing the edge between X and Y from Y to X will inevitably produce an additional collider

$$V \rightarrow X \leftarrow Z$$

or a directed cycle.

Hence, the edge must be directed from X to Y .

STEP 4: IDENTIFY DERIVED DIRECTIONS.

--- X and V are adjacent nodes,
i.e., $X \rightarrow V$, $V \rightarrow X$ or $X - V$.

It follows from the fact that directing the edge between X and Y from Y to X will inevitably produce an additional collider

$Y \rightarrow X \leftarrow Z$

or a directed cycle.

Hence, the edge must be directed from X to Y .

The dashed lines used to illustrate the fourth rule indicate that X and V are connected by an edge (either directed or not, $X \rightarrow V$, $V \rightarrow X$ or $X - V$).

The fourth rule is not necessary if the orientation of the initial PDAG is limited to containing colliders only.

The initial PDAG may contain non-colliders when expert knowledge on edge directions are included in the graph.

STEP 4: IDENTIFY DERIVED DIRECTIONS.

As neither the collider

$$B \rightarrow A \leftarrow W$$

nor the collider

$$E \rightarrow A \leftarrow W$$

was identified as a collider of \mathcal{G} , the edge between A and W must be directed from A to W (Figure 12.13). (Application of **RULE 1**)

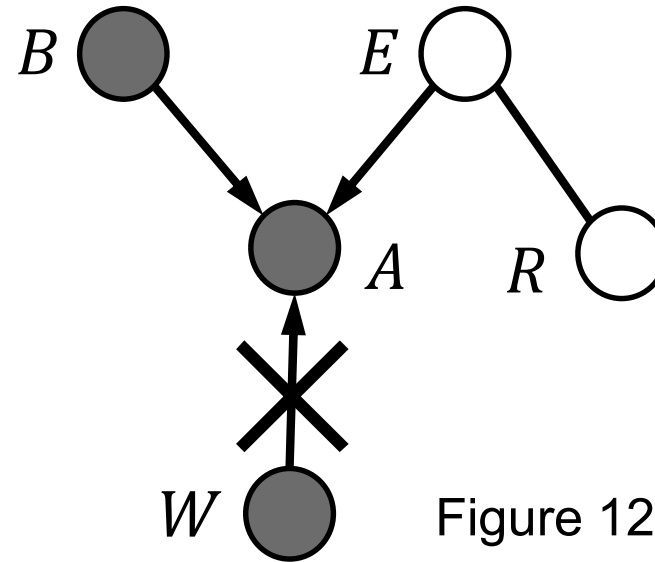


Figure 12.12

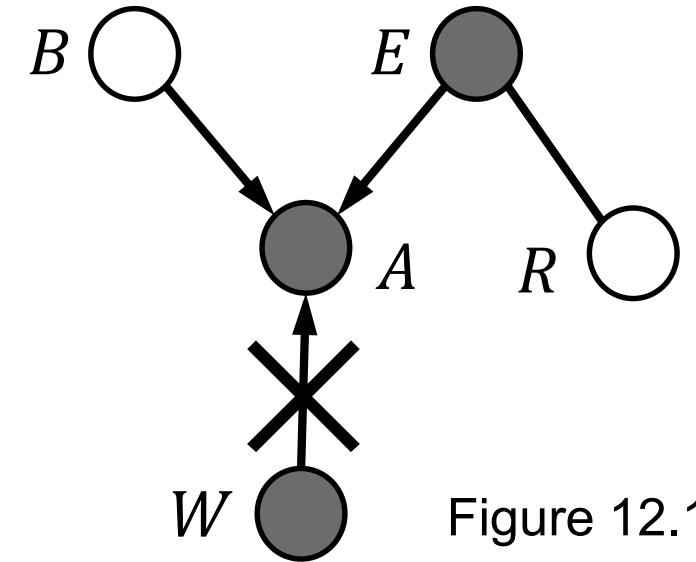


Figure 12.15

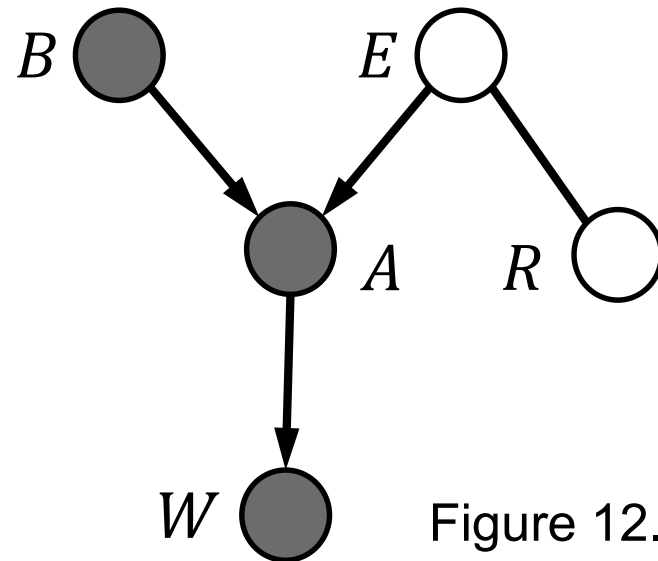


Figure 12.13

STEP 4: IDENTIFY DERIVED DIRECTIONS.

As neither the collider

$$B \rightarrow A \leftarrow W$$

nor the collider

$$E \rightarrow A \leftarrow W$$

was identified as a collider of \mathcal{G} , the edge between A and W must be directed from A to W (Figure 12.14). (Application of **RULE 1**)

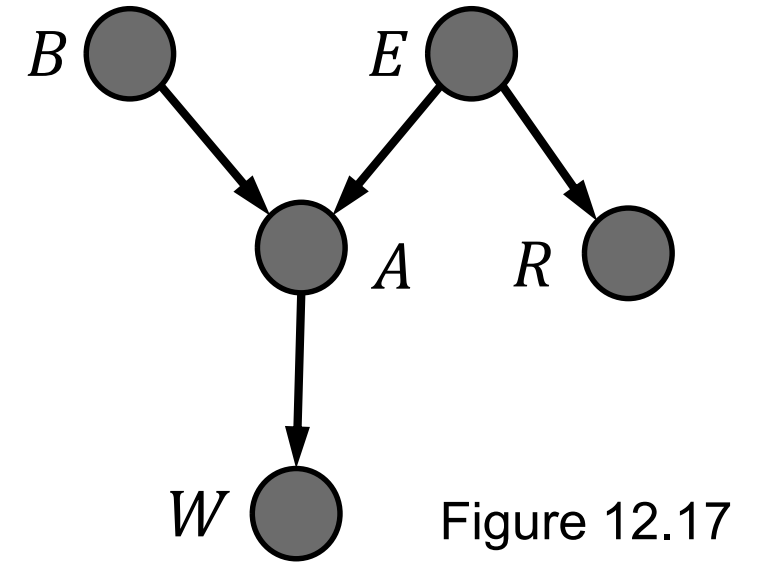
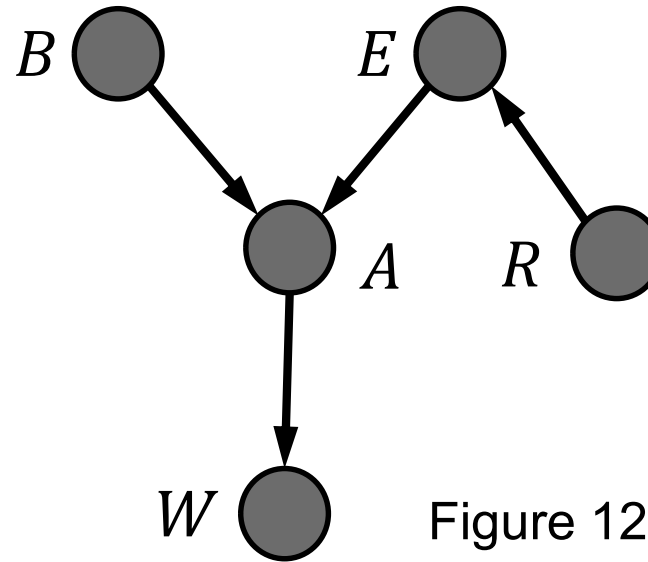
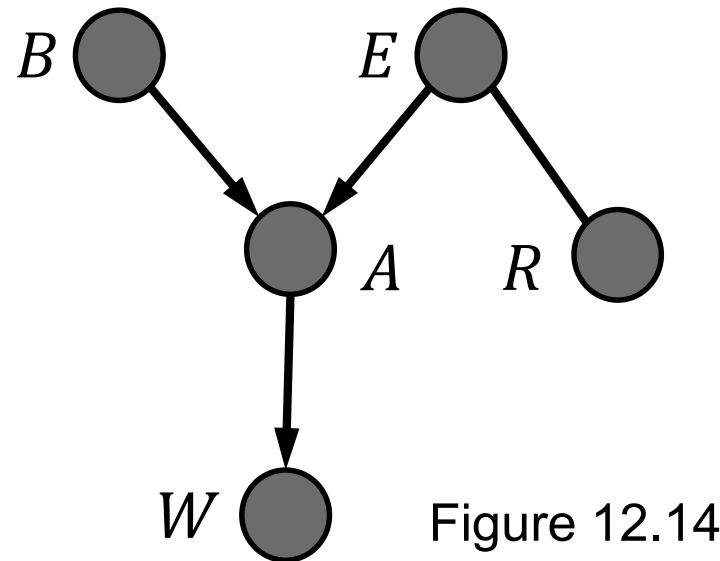
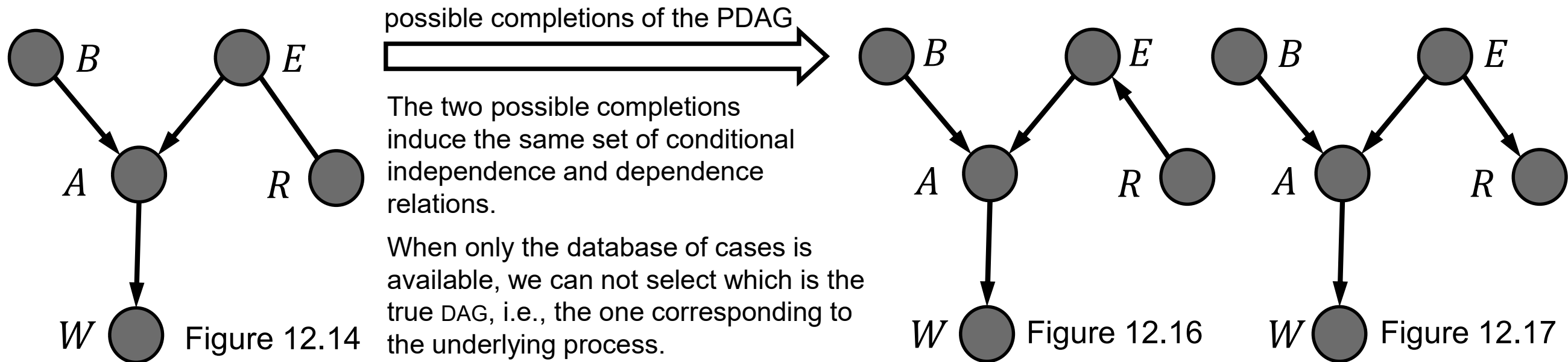


Figure 12.16 and Figure 12.17 show the **EQUIVALENCE CLASS** of $\mathcal{M}_{\mathcal{D}}$.

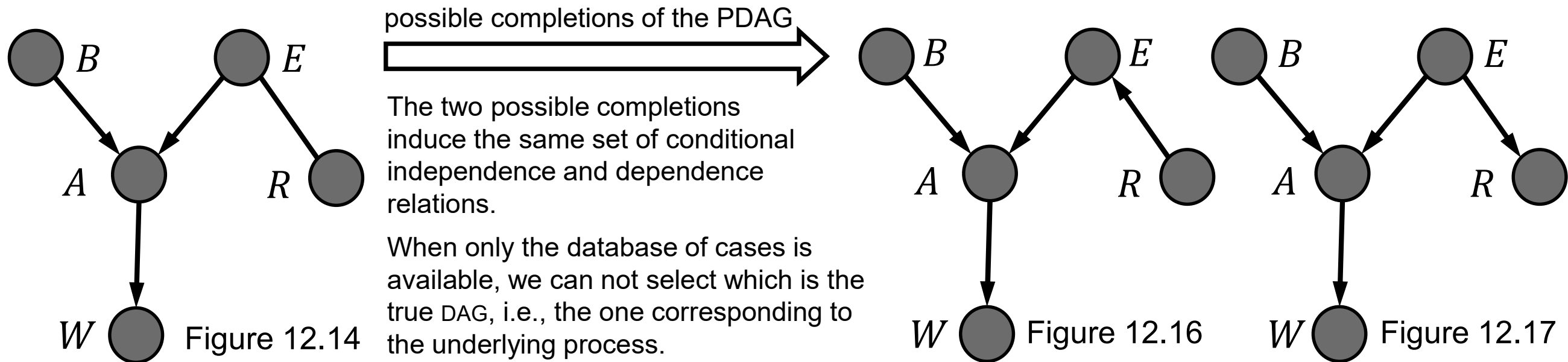
The **EQUIVALENCE CLASS** of $\mathcal{M}_{\mathcal{D}}$ contains two DAGs differing only with respect to the orientation of the edge between E and R .

- The four rules are necessary and sufficient for achieving maximal orientation (up to equivalence) of the PDAG returned by the PC algorithm.
- We use these four rules repeatedly until no edge can be given an orientation.
- Notice that the result of closing edge directions under rules from 1 to 4 is not necessarily a DAG.
- If the graph is not a DAG, then expert knowledge may be appropriate in order to direct an edge.
- Once an edge has been directed by use of expert knowledge derived directions should be identified.
- This process may be repeated until a DAG structure is obtained.
- Experience shows that most edges are directed using **RULE 1**, and that **RULE 3** is only rarely used.



- When a unique DAG can not be obtained from the database of cases we can ask help to the domain expert.
- If we can obtain a unique DAG, then we move on to learn the parameters of the causal network.

- There are **OTHER CONSTRAINT-BASED ALGORITHMS** that allow us to drop various assumptions.
- The FCI (**FAST CAUSAL INFERENCE**) algorithm works without assuming causal sufficiency.
- The **CCD** algorithm works without assuming acyclicity.
- **CONSTRAINT-BASED ALGORITHMS** suffer from the fact that conditional independence tests are hard, and it can sometimes require a lot of data to get accurate test results.



Progress in Artificial Intelligence
<https://doi.org/10.1007/s13748-019-00194-y>

REVIEW



A survey on Bayesian network structure learning from data

Mauro Scanagatta¹  · Antonio Salmerón² · Fabio Stella³

Received: 14 March 2019 / Accepted: 20 May 2019
© Springer-Verlag GmbH Germany, part of Springer Nature 2019



Abstract

A necessary step in the development of artificial intelligence is to enable a machine to *represent* how the world works, building an internal structure from data. This structure should hold a good trade-off between expressive power and querying efficiency. Bayesian networks have proven to be an effective and versatile tool for the task at hand. They have been applied to modeling knowledge in a variety of fields, ranging from bioinformatics to law, from image processing to economic risk analysis. A crucial aspect is learning the dependency graph of a Bayesian network from data. This task, called *structure learning*, is NP-hard and is the subject of intense, cutting-edge research. In short, it can be thought of as choosing one graph over the many candidates, grounding our reasoning over a collection of samples of the distribution generating the data. The number of possible graphs increases very quickly at the increase in the number of variables. Searching in this space, and selecting a graph over the others, becomes quickly burdensome. In this survey, we review the most relevant structure learning algorithms that have been proposed in the literature. We classify them according to the approach they follow for solving the problem and we also show alternatives for handling missing data and continuous variable. An extensive review of existing software tools is also given.

Keywords Machine learning · Statistics · Bayesian network · Structure learning