

# Advanced Probabilistic Modelling

## Definitions and Fundamentals

---

Marco Scutari

Dalle Molle Institute for  
Artificial Intelligence (IDSIA)

**Machine learning** studies the algorithms and the statistical tools that allow computer systems to perform specific, well-defined tasks without explicit instructions. It is a sub-field of **artificial intelligence**.

Broadly speaking, in order to do this:

1. We need a **working model of the world** that describes the task and its context in a way a computer can understand.
2. We need a goal: how do we **measure the performance** of the model? Because that is what we optimise for! Usually it is the ability to **predict new events**.
3. We **encode our knowledge** of the world drawing information from training data, experts or both: this is called **learning**.
4. The computer system uses the model as a **proxy of reality** and to **perform inference** and decide if/how to perform the assigned task as new inputs come in.

## IDENTIFY THE VARIABLES TO INCLUDE IN THE MODEL

---

The first step in building a machine learning model is to choose which variables to include. **Which aspects of/entities in the world do we need the model to represent** for the computer to carry out the assigned task? This is known as **feature selection**.

- **Each** aspect of the world or entity is modelled with **one random variable**.
- We should use a **small enough number** of variables because if we have too many:
  - it is difficult it is to construct the model;
  - it is difficult to interpret and to troubleshoot it;
  - the model requires too much computing power to learn and to run.
- We must choose which are the **relevant events** that make up the sample space of each variable, again taking care of not having too many.

## AN EXAMPLE: THE CAR START PROBLEM

---

Fuel

Spark Plugs

Fuel Meter

Start

	Realistic	Pragmatic
Fuel	0%–100%	Yes, No
Spark Plugs	Work, Fault	Work, Fault
Start	Yes, No	Yes, No
Fuel Meter	0%–100%	Empty, Half, Full

The second step is choosing which class of machine learning models to select from.

- **Generative models:** we have a set of variables  $X_1, \dots, X_N$  describing various components of a complex phenomenon, and we are interested in building a mechanistic model of that phenomenon to **understand it**. Therefore, we want to show how the various parts interact with each other. In order to do so we choose to model their joint probability  $P(X_1, \dots, X_N)$ .
- **Discriminative models:** we have one particular variable (say,  $X_1$ ) that is closely tied with our model task, and a number of other variables ( $X_2, \dots, X_N$ ) which we believe can be used to **predict it**. We do not care about how the  $X_i$  are related to each other, so we just model  $P(X_1 | X_2, \dots, X_N)$ .

How do we decide whether there is a relationships between variables?

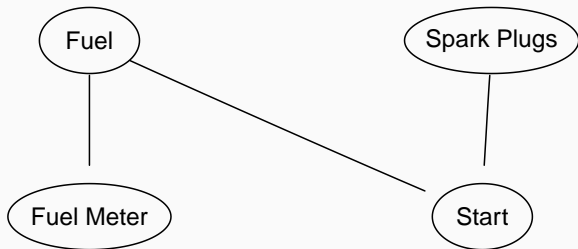
We never have perfect knowledge of what we are modelling: hence we use the language of **probability**, and we say that two variables are associated if the occurrence of an event in one variable affects the probability of an event occurring in another variable, possibly conditional on other variables.

How can we acquire information on what we are modelling:

- consulting domain **experts**;
- using probability and statistics to extract it from **data**;
- a combination of **both**.

## THE CAR START PROBLEM, WITH EDGES

---



- The Fuel Meter measures the amount of Fuel?
- The Spark Plugs ignite the Fuel to Start the car?
- If there is Fuel in the car, it can start even if the Fuel Meter is wrong and displays 0%?

In probability **associations are symmetric**: the derivation of Bayes' theorem makes it really clear that

$$P(X_1 | X_2) P(X_2) = P(X_1, X_2) = P(X_2 | X_1) P(X_1).$$

However, it feels more natural to choose the conditioning variables such that they affect the conditioned variables instead of the other way round.

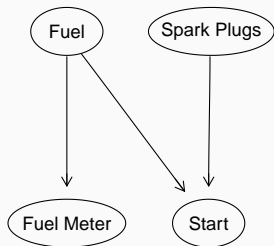
But what does that mean from a modelling point of view? It means that we are **giving arcs a causal interpretation** and that we **choose arc directions to go from cause (nodes) to effect (nodes)**.

How do we do that?

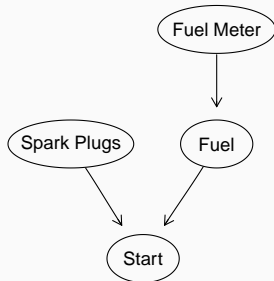


## THE CAR START PROBLEM, WITH ARCS

$$\begin{aligned} P(\text{Start}, \text{Fuel Meter}, \text{Fuel}, \text{Spark Plugs}) \\ = P(\text{Start} | \text{Fuel}, \text{Spark Plugs}) \times \\ P(\text{Fuel Meter} | \text{Fuel}) P(\text{Fuel}) \times \\ P(\text{Spark Plugs}) \end{aligned}$$



$$\begin{aligned} P(\text{Start}, \text{Fuel Meter}, \text{Fuel}, \text{Spark Plugs}) \\ = P(\text{Start} | \text{Fuel}, \text{Spark Plugs}) \times \\ P(\text{Fuel Meter}) P(\text{Fuel} | \text{Fuel Meter}) \times \\ P(\text{Spark Plugs}) \end{aligned}$$



The criterion to identify causes and effect is **intervention**. Consider:

- If we fill the tank with fuel, the fuel meter goes up.
- If we tamper with the fuel meter to make it say Full, the fuel tank does not magically refill itself.

Hence, Fuel is the **cause** and Fuel Meter is the effect and the most intuitive arc direction is Fuel  $\rightarrow$  Fuel Meter.

What the probability  $P(\text{Fuel Meter} \mid \text{Fuel})$  tells us is just that if the fuel meter says Full there probably is fuel in the tank, whereas if the fuel meter says Empty there may be no fuel in the tank (assuming the fuel meter works reliably).

## CAR START: THE CONDITIONAL PROBABILITIES

Spark Plugs		Fuel		Start	
Work	Fault	Yes	No	Spark Plugs = Work	
				Fuel = Yes	Fuel = No
?	?	?	?	Yes	?
				No	?

Fuel Meter		Start		
	Fuel = Yes	Fuel = No	Spark Plugs = Fault	
			Fuel = Yes	Fuel = No
Empty	?	?	Yes	?
Half	?	?	No	?
Full	?	?		?

After we decide that the first model is good to go, we need to:

1. choose **which distribution to use for each node**;
2. fill in the values of its **parameters** by asking domain experts, estimating them from data or a combination of the two.

The number of parameters gives the **complexity** of the model, rather than the number of nodes or the number of arcs.

A more general way of using a model is to **interrogate** it: we **have some evidence** on some of the variables (that is, we assume we know their values), and we would like to know the **the probability of some event**.

For instance: say that Fuel Meter = Half. How does  $P(\text{Start} = \text{Yes})$  change after we **introduce this evidence in the model**?

Predicting Start from all the other variables is a particular case in which we have evidence on all the other variables.

## CAR START: THE EXHAUSTIVE (DUMB) WAY

---

Armed with patience, we start by writing

$$P(\text{Start} = \text{Yes}) = \\ P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}) + P(\text{Start} = \text{Yes}, \text{Fuel} = \text{No})$$

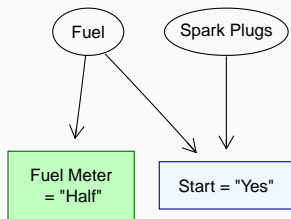
and then, **recursively**,

$$P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}) = \\ P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}, \text{Spark.Plugs} = \text{Work}) + \\ P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}, \text{Spark.Plugs} = \text{Fault})$$

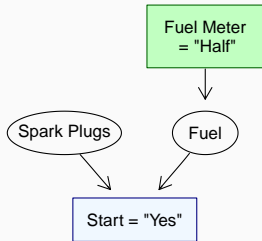
$$P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}, \text{Spark.Plugs} = \text{Work}) \\ = P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}, \text{Spark.Plugs} = \text{Work}, \text{Fuel.Meter} = \text{Full}) + \\ P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}, \text{Spark.Plugs} = \text{Work}, \text{Fuel.Meter} = \text{Half}) + \\ P(\text{Start} = \text{Yes}, \text{Fuel} = \text{Yes}, \text{Spark.Plugs} = \text{Work}, \text{Fuel.Meter} = \text{Empty})$$

## CAR START: THE PRINCIPLED (PROBABILISTIC) WAY

$$\begin{aligned} &P(\text{Start} = \text{Yes}, \text{Fuel Meter} = \text{Half}, \\ &\quad \text{Fuel}, \text{Spark Plugs}) = \\ &= P(\text{Start} = \text{Yes} \mid \text{Fuel}, \text{Spark Plugs}) \times \\ &\quad P(\text{Fuel Meter} = \text{Half} \mid \text{Fuel}) \times \\ &\quad P(\text{Fuel}) P(\text{Spark Plugs}) \end{aligned}$$

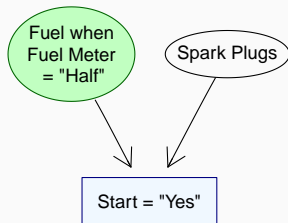


$$\begin{aligned} &P(\text{Start} = \text{Yes}, \text{Fuel Meter} = \text{Half}, \\ &\quad \text{Fuel}, \text{Spark Plugs}) \\ &= P(\text{Start} \mid \text{Fuel}, \text{Spark Plugs}) \times \\ &\quad P(\text{Fuel} \mid \text{Fuel Meter} = \text{Half}) \times \\ &\quad \frac{P(\text{Fuel Meter} = \text{Half})}{\cancel{P(\text{Fuel})}} \times \\ &\quad \cancel{P(\text{Fuel})} P(\text{Spark Plugs}) \end{aligned}$$



## CAR START: THE PRINCIPLED (PROBABILISTIC) WAY

$$\begin{aligned} &P(\text{Start} = \text{Yes}, \text{Fuel}, \text{Spark Plugs} | \\ &\quad \text{Fuel Meter} = \text{Half}) \\ &= P(\text{Start} | \text{Fuel}, \text{Spark Plugs}) \times \\ &\quad P(\text{Fuel} | \text{Fuel Meter} = \text{Half}) \times \\ &\quad \frac{P(\text{Fuel Meter} = \text{Half})}{P(\text{Fuel Meter} = \text{Half})} \times \\ &\quad P(\text{Spark Plugs}) \end{aligned}$$



This leaves three variables, of which Start is fixed to Yes: hence we have to consider  $P(\text{Start} = \text{Yes})$  under **four scenarios**:

Fuel = Yes | Fuel Meter = Half, Spark Plugs = Work

Fuel = Yes | Fuel Meter = Half, Spark Plugs = Fault

Fuel = No | Fuel Meter = Half, Spark Plugs = Work

Fuel = No | Fuel Meter = Half, Spark Plugs = Fault

and **sum the corresponding  $P(\text{Start} = \text{Yes} | \text{scenario}) P(\text{scenario})$ .**

Exhaustive enumeration obviously:

1. does not scale (try that with 20 variables!);
2. is only feasible in the first place if all variables are discrete.

Each of the steps in the previous slide corresponds to both

- an (probabilistic) application of Bayes theorem
- a (graphical) manipulation of arcs and nodes.

We can use the fact that arcs represent probabilistic associations to perform symbolic computations through graphical operations!

We can also treat this model like a hierarchical model and adapt the literature on **Monte Carlo simulations**.



Now we want to automate the whole process, so that the computer system itself will (ideally) do all the work.

A model that promises to do this is Bayesian networks (BNs):

- They combine graphs and probability as we did earlier, but in a rigorous fashion.
- There are automated reasoning algorithms for that use the graphical part of the model to guide a computer system in manipulating probability distributions, computing probabilities of and predicting events of interest.
- It is possible to learn them automatically from data.
- They can be used as causal models.
- As far as models, go they are very green: they recycle large amounts of results from classical statistics.

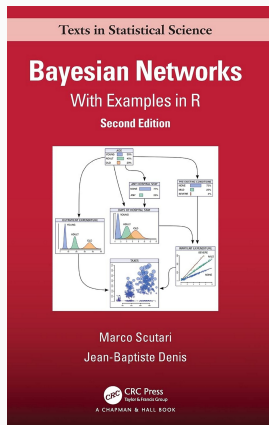
The most comprehensive R package for working with Bayesian networks is **bnlearn**, which you should install by

```
install.packages("bnlearn")
```

The reference **website** for **bnlearn** is:

<http://www.bnlearn.com>

And there is a reference **book** too!



Bayesian networks (BNs) are defined by:

- a **network structure**, a **directed acyclic graph**  $\mathcal{G}$ , in which each node corresponds to a random variable  $X_i$ ;
- a **global probability distribution** over  $\mathbf{X} = \{X_1, \dots, X_N\}$  which can be factorised into smaller **local probability distributions** according to the arcs present in the graph.

The main role of the network structure is to express the **conditional independence** relationships among the variables in the model through **graphical separation**, thus specifying the factorisation of the global distribution:

$$P(\mathbf{X}; \Theta) = \prod_{i=1}^N P(X_i | \Pi_{X_i}; \Theta_{X_i}) \quad \text{where} \quad \Pi_{X_i} = \{\text{parents of } X_i\}.$$

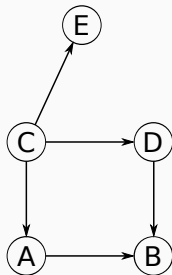
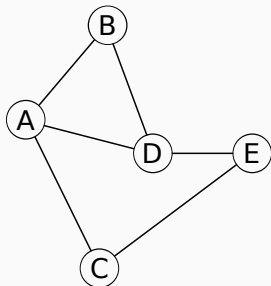
The first component of a BN is a graph. A graph  $\mathcal{G}$  is a mathematical object with:

- a set of **nodes**;
- a set of **arcs**  $A$  which are identified by pairs for nodes.

Given the nodes, a graph is uniquely identified by the arc set. An arc can be:

- **undirected** if the arc has no direction, for instance  $A - B$ ;
- **directed** if the arc has a specific direction, for instance  $A \rightarrow B$ .

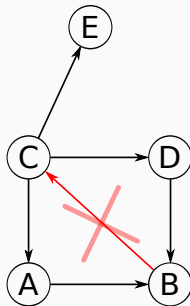
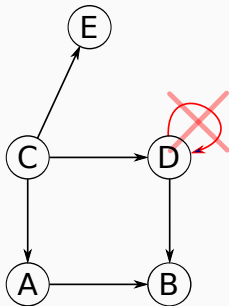
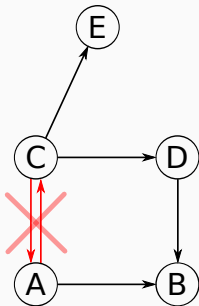
The assumption is that there is at most one arc between each pair of nodes.



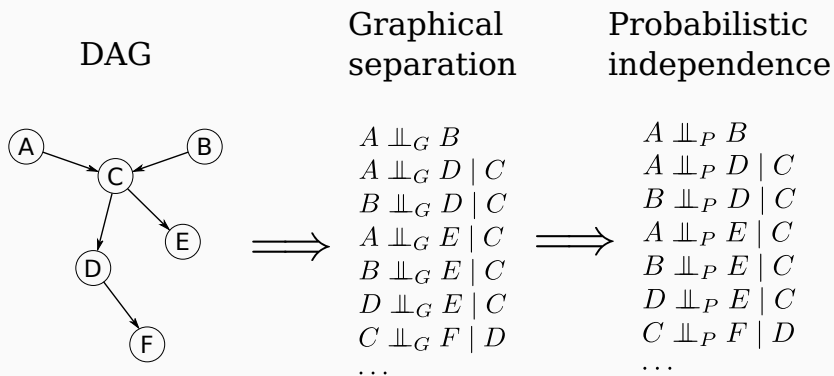
## DIRECTED ACYCLIC GRAPHS

BNs use a specific kind of graph called a **directed acyclic graph** (DAG), that:

- contains only directed arcs;
- does not contain any loop (an arc  $D \rightarrow D$  from a node to itself);
- does not contain any cycle (a sequence of arcs like  $B \rightarrow C \rightarrow D \rightarrow B$  that starts and ends in the same node).

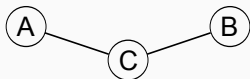


## HOW THE DAG MAPS TO THE PROBABILITY DISTRIBUTION



Formally, the DAG is an **independence map** of the probability distribution of  $\mathbf{X}$ , with graphical separation ( $\perp\!\!\!\perp_G$ ) implying probabilistic independence ( $\perp\!\!\!\perp_P$ ).

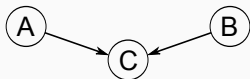
separation (undirected graphs)



$$A \perp\!\!\!\perp B \mid C$$

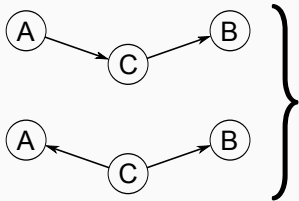
$$P(A, B, C) = P(A \mid C) P(B \mid C) P(C)$$

d-separation (directed acyclic graphs)



$$A \not\perp\!\!\!\perp B \mid C$$

$$P(A, B, C) = P(C \mid A, B) P(A) P(B)$$



$$A \perp\!\!\!\perp B \mid C$$

$$\begin{aligned} P(A, B, C) &= \\ &= P(B \mid C) P(C \mid A) P(A) \\ &= P(A \mid C) P(B \mid C) P(C) \end{aligned}$$

Now, in the general case we can extend the patterns from the fundamental connections and apply them to every possible path between **A** and **B** for a given **C**; this is how **d-separation** is defined.

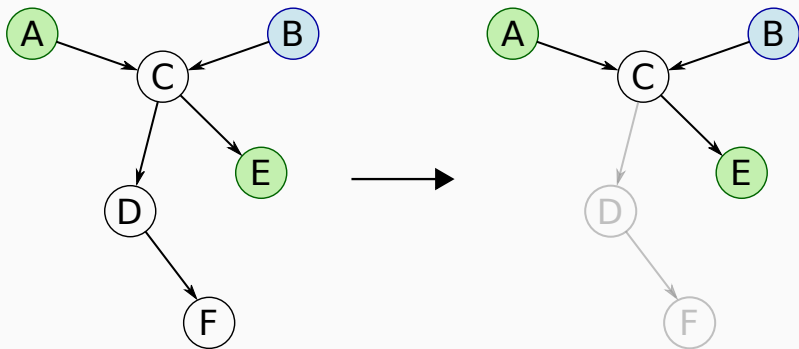
*If **A**, **B** and **C** are three disjoint subsets of nodes in a directed acyclic graph  $\mathcal{G}$ , then **C** is said to d-separate **A** from **B**, denoted  $\mathbf{A} \perp\!\!\!\perp_{\mathcal{G}} \mathbf{B} \mid \mathbf{C}$ , if along every path between a node in **A** and a node in **B** there is a node  $v$  satisfying one of the following two conditions:*

- 1.  $v$  has converging edges (that is, there are two edges pointing to  $v$  from the adjacent nodes in the path) and none of  $v$  or its descendants (that is, the nodes that can be reached from  $v$ ) are in **C**.*
- 2.  $v$  is in **C** and does not have converging edges.*

This definition clearly **does not provide a computationally feasible approach** to assess d-separation; but there are other ways.

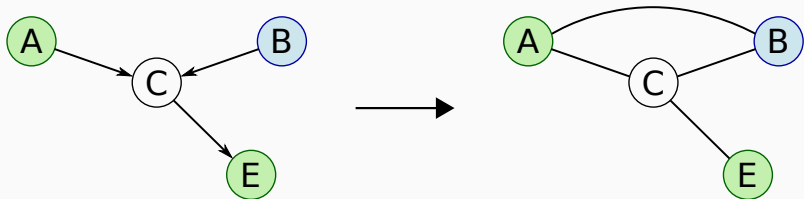


## A SIMPLE ALGORITHM TO CHECK D-SEPARATION



Say that we want to check whether  $A$  and  $E$  are d-separated by  $B$ . First, we can **drop all the nodes that are not ancestors** (that is, parents, parents' parents, etc.) of  $A$ ,  $E$  and  $B$  since each node only depends on its parents.

## A SIMPLE ALGORITHM TO CHECK D-SEPARATION



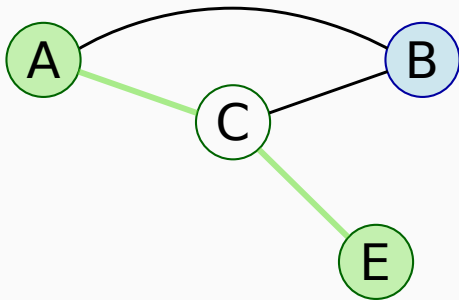
We then transform the subgraph into its **moral graph** by

1. connecting all the nodes that have one child in common; and
2. removing all arc directions to obtain an undirected graph.

This transformation makes the dependence between parents explicit by “marrying” them and of makes it possible for us to use the classic definition of graphical separation.

## A SIMPLE ALGORITHM TO CHECK D-SEPARATION

---



Finally, we can just perform a depth-first or breadth-first search and see **if we can find an open path** between  $A$  and  $E$ , that is, a path that is not blocked by  $B$ .

## D-SEPARATION EXAMPLE: THE DAG WE CREATED EARLIER

---

The last graph is an undirected graph: if there is a path from A to E there is a path from E to A. This means that **d-separation is symmetric**:

$$A \not\perp_G E \mid B \iff E \not\perp_G A \mid B$$

Which must be the case because **independence is also symmetric**,

$$P(A, E \mid B) = P(E, A \mid B) \neq P(A \mid B) P(E \mid B),$$

and d-separation implies probabilistic independence.

**NOTE:** d-separation does not necessarily require a separating set. Or, to put it in another way, the separating set can be empty. In that case we are checking whether variables are **marginally independent** because there is no path at all that connects them.

If we use d-separation as our definition of graphical separation, assuming that the DAG is an independence map leads to the general formulation of the **decomposition of the global distribution**

$$P(\mathbf{X}) = \prod_{i=1}^N P(X_i \mid \Pi_{X_i})$$

into the **local distributions** for the  $X_i$  given their parents  $\Pi_{X_i}$ . If  $X_i$  has two or more parents it depends on their joint distribution, because each pair of parents forms a convergent connection centred on  $X_i$  and we cannot establish their independence. This decomposition is preferable to that obtained from the chain rule,

$$P(\mathbf{X}) = \prod_{i=1}^N P(X_i \mid X_{i+1}, \dots, X_N),$$

because the conditioning sets are typically smaller.

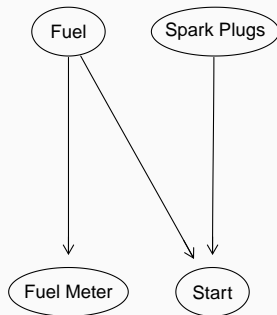
Another result along the same lines is the **local Markov property**, which can be combined with the chain rule above to get the decomposition into local distributions.

*Each node  $X_i$  is conditionally independent of its non-descendants (the nodes  $X_j$  for which there is no path from  $X_i$  to  $X_j$ ) given its parents.*

Compared to the previous decomposition, it highlights the fact that parents are not completely independent from their children in the BN: a trivial application of Bayes' theorem to invert the direction of the conditioning shows how information on a child can change the distribution of the parent.

## THE LOCAL MARKOV PROPERTY: CAR START

---



The **parent sets**:

$$\text{Fuel} = \{\}$$

$$\text{Fuel Meter} = \{\text{Fuel}\}$$

$$\text{Spark Plugs} = \{\}$$

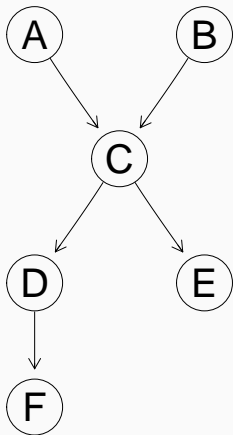
$$\text{Start} = \{\text{Fuel}, \text{Spark Plugs}\}$$

The corresponding **decomposition**:

$$\begin{aligned} P(\text{Start}, \text{Fuel Meter}, \text{Fuel}, \\ \text{Spark Plugs}) = \\ P(\text{Start} | \text{Fuel}, \text{Spark Plugs}) \\ P(\text{Fuel Meter} | \text{Fuel}) \times \\ P(\text{Fuel}) P(\text{Spark Plugs}) \end{aligned}$$

## THE LOCAL MARKOV PROPERTY: THE DAG WE CREATED EARLIER

---



The **parent sets**:

$$A = \{\}$$

$$B = \{\}$$

$$C = \{A, B\}$$

$$D = \{C\}$$

$$E = \{C\}$$

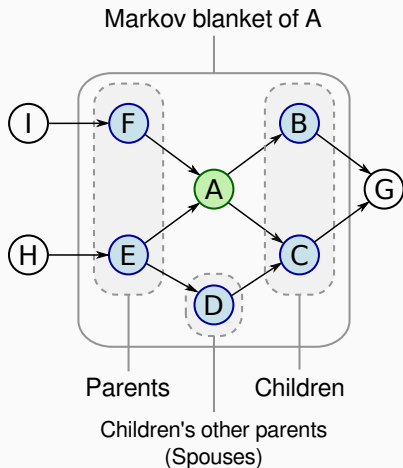
$$F = \{D\}$$

The corresponding **decomposition**:

$$\begin{aligned} P(A, B, C, D, E, F) = & P(A) P(B) P(C | A, B) \\ & P(D | C) P(E | C) P(F | D) \end{aligned}$$



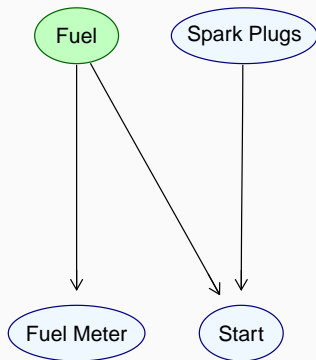
## COMPLETELY D-SEPARATING: MARKOV BLANKETS



We can easily use the DAG to solve the **feature selection** problem. The set of nodes that graphically isolates a target node from the rest of the DAG is called its **Markov blanket** and includes:

- its parents;
- its children;
- other nodes sharing a child.

Since  $\perp\!\!\!\perp_G$  implies  $\perp\!\!\!\perp_P$ , we can restrict ourselves to the Markov blanket to perform any kind of inference on the target node, and disregard the rest.



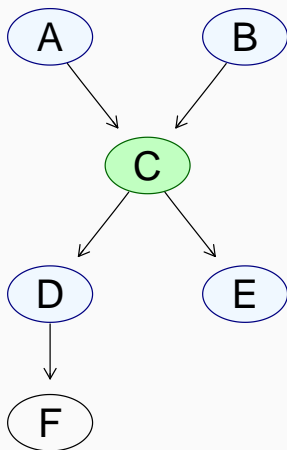
The **parents, children** and **spouses** of Fuel:

$\{\}$   
 $\{\text{Fuel Meter, Start}\}$   
 $\{\text{Spark Plugs}\}$

The **Markov blanket** of Fuel:  
 $\{\text{Fuel Meter, Spark Plugs, Start}\}$

## MARKOV BLANKET: THE DAG WE CREATED EARLIER

---



Printing the **parents**, **children**  
and **spouses** of C:

$\{A, B\}$

$\{D, E\}$

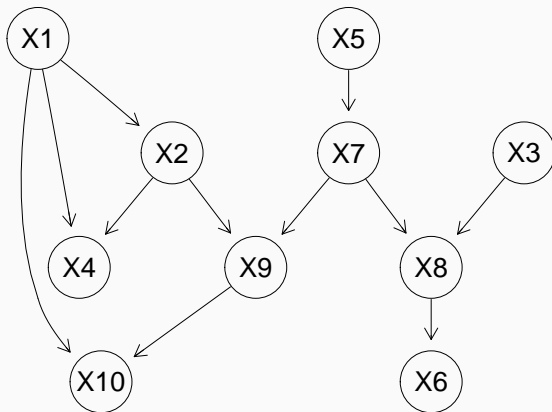
$\{\}$

The **Markov blanket** of C:

$\{A, B, D, E\}$

## DIFFERENT DAGs, SAME DISTRIBUTION

A DAG uniquely identifies a factorisation of  $P(\mathbf{X})$ ; **the converse is not necessarily true**. Consider this DAG:



The decomposition into local distributions is:

$$\begin{aligned} P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9, X_{10}) = \\ P(X_1) P(X_3) \underbrace{P(X_5)}_{X_5} P(X_6 | X_8) P(X_2 | X_1) \underbrace{P(X_7 | X_5)}_{X_5 \rightarrow X_7} \\ P(X_4 | X_1, X_2) P(X_8 | X_3, X_7) P(X_9 | X_2, X_7) P(X_{10} | X_1, X_9). \end{aligned}$$

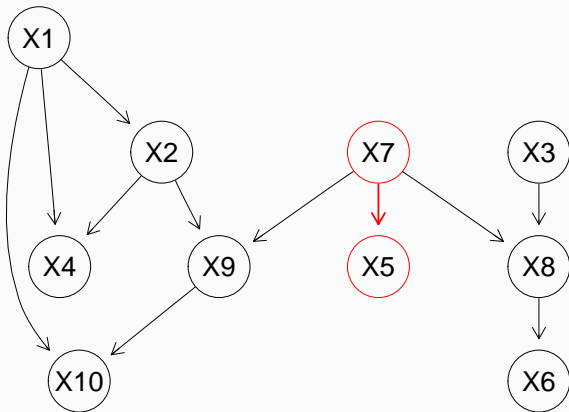
However, **look at  $X_5 \rightarrow X_7$** :  $P(X_7 | X_5) P(X_5) = P(X_5 | X_7) P(X_7)$  by Bayes' theorem. Then

$$\begin{aligned} P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9, X_{10}) = \\ P(X_1) P(X_3) \underbrace{P(X_7)}_{X_7} P(X_6 | X_8) P(X_2 | X_1) \underbrace{P(X_5 | X_7)}_{X_7 \rightarrow X_5} \\ P(X_4 | X_1, X_2) P(X_8 | X_3, X_7) P(X_9 | X_2, X_7) P(X_{10} | X_1, X_9). \end{aligned}$$

## DIFFERENT DAGs, SAME DISTRIBUTION

---

The DAG that gives this **new, equivalent decomposition** is:



Next let's look at  $X_8 \rightarrow X_6$ .

$$\begin{aligned} P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9, X_{10}) = \\ P(X_1) P(X_3) P(X_7) \underbrace{P(X_6 | X_8)}_{X_8 \rightarrow X_6} P(X_2 | X_1) P(X_5 | X_7) \\ P(X_4 | X_1, X_2) \underbrace{P(X_8 | X_3, X_7)}_{X_8 \leftarrow X_3, X_8 \leftarrow X_7} P(X_9 | X_2, X_7) P(X_{10} | X_1, X_9). \end{aligned}$$

We cannot reverse the  $X_8 \rightarrow X_6$  as we did with  $X_5 \rightarrow X_7$  without changing the probability distribution. If we try, we get

$$P(X_6 | X_8) P(X_8 | X_3, X_7) = P(X_8 | X_6) P(X_6) \frac{P(X_8 | X_3, X_7)}{P(X_8)},$$

which does not simplify because  $X_8$  has other parents ( $X_3, X_7$ ).

Finally, let's look at  $X_1$ ,  $X_2$  and  $X_4$ .

$$\begin{aligned} P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9, X_{10}) = \\ \underbrace{P(X_1)}_{X_1} P(X_3) P(X_5) P(X_6 | X_8) \underbrace{P(X_2 | X_1)}_{X_1 \rightarrow X_2} P(X_7 | X_5) \\ \underbrace{P(X_4 | X_1, X_2)}_{X_1 \rightarrow X_4, X_2 \rightarrow X_4} P(X_8 | X_3, X_7) P(X_9 | X_2, X_7) P(X_{10} | X_1, X_9). \end{aligned}$$

By Bayes' theorem we can say

$$\begin{aligned} P(X_1) P(X_2 | X_1) P(X_4 | X_1, X_2) = P(X_1, X_2, X_4) = \\ \underbrace{P(X_4)}_{X_4} \underbrace{P(X_2 | X_4)}_{X_4 \rightarrow X_2} \underbrace{P(X_1 | X_2, X_4)}_{X_1 \leftarrow X_2, X_1 \leftarrow X_4} \end{aligned}$$

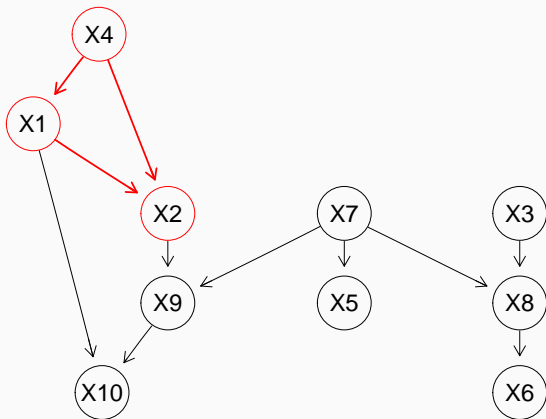
which gives us another DAG again.



## DIFFERENT DAGs, SAME DISTRIBUTION

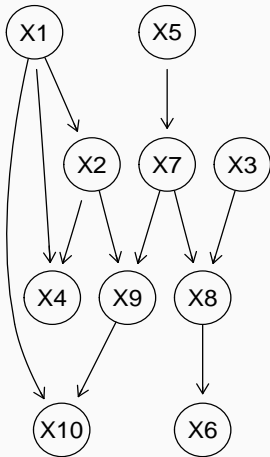
---

The DAG that gives this **last equivalent decomposition** is:

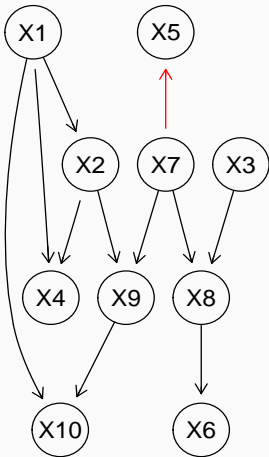


## COMPARING THESE DIFFERENT DAGs

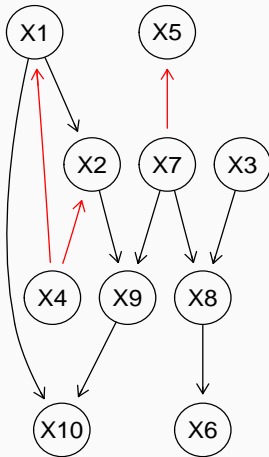
original



equivalent #1



equivalent #2



## DIFFERENT DAGs, SAME DISTRIBUTION: EQUIVALENCE CLASSES

---

To sum it up: we can reverse a number of arcs without changing the dependence structure of  $\mathbf{X}$ . Since the fundamental connections  $A \rightarrow C \rightarrow B$  and  $A \leftarrow C \rightarrow B$  are probabilistically equivalent, we can reverse the directions of their arcs as we like as long as we do not create any new **v-structure** ( $A \rightarrow C \leftarrow B$ , with no arc between  $A$  and  $B$ ).

This means that we can group DAGs into **equivalence classes** that are uniquely identified by the underlying undirected graph and the v-structures. The directions of other arcs can be either:

- uniquely identifiable because one of the directions would introduce cycles or new v-structures (**compelled arcs**);
- completely undetermined.

The result is a **completed partially directed graph** (CPDAG).

## WHAT ARE V-STRUCTURES, AND WHAT ARE NOT

---

It is important to note that even though  $A \rightarrow C \leftarrow B$  is a convergent connection, **it is not a v-structure if  $A$  and  $C$  are connected by  $A \rightarrow B$  or  $B \rightarrow A$** . In that case, we are no longer able to identify which nodes are the parents in the connection.

For instance:

$$\begin{aligned} \underbrace{P(A) P(B|A) P(C|A, B)}_{A \rightarrow C \leftarrow B, A \rightarrow B} &= P(A) \frac{P(B, A)}{P(A)} \frac{P(C, A, B)}{P(A, B)} = \\ &= P(A) P(C, B|A) = \underbrace{P(A) P(B|C, A) P(C|A)}_{C \rightarrow B \leftarrow A, A \rightarrow C}. \end{aligned}$$

Therefore, the fact that the two parents in a v-structure are not connected is crucial in the identification of the correct CPDAG.

From this description we can tell different groups of arcs apart:

Directed arcs:

$X_8 \rightarrow X_6$   
 $X_1 \rightarrow X_2$   
 $X_5 \rightarrow X_7$   
 $X_1 \rightarrow X_4$   
 $X_2 \rightarrow X_4$   
 $X_3 \rightarrow X_8$   
 $X_7 \rightarrow X_8$   
 $X_2 \rightarrow X_9$   
 $X_7 \rightarrow X_9$   
 $X_1 \rightarrow X_{10}$   
 $X_9 \rightarrow X_{10}$

Undirected arcs:

None.

Compelled arcs:

$X_1 \rightarrow X_{10}$   
 $X_2 \rightarrow X_9$   
 $X_3 \rightarrow X_8$   
 $X_7 \rightarrow X_8$   
 $X_7 \rightarrow X_9$   
 $X_8 \rightarrow X_6$   
 $X_9 \rightarrow X_{10}$

V-structures:

$X_1 \rightarrow X_{10} \leftarrow X_9$   
 $X_3 \rightarrow X_8 \leftarrow X_7$   
 $X_2 \rightarrow X_9 \leftarrow X_7$

## THE CORRESPONDING CPDAG

---

Which in the corresponding CPDAG become:

Directed arcs:

$X_1 \rightarrow X_{10}$   
 $X_2 \rightarrow X_9$   
 $X_3 \rightarrow X_8$   
 $X_7 \rightarrow X_8$   
 $X_7 \rightarrow X_9$   
 $X_8 \rightarrow X_6$   
 $X_9 \rightarrow X_{10}$

Undirected arcs:

$X_1 - X_2$   
 $X_1 - X_4$   
 $X_2 - X_1$   
 $X_2 - X_4$   
 $X_4 - X_1$   
 $X_4 - X_2$   
 $X_5 - X_7$   
 $X_7 - X_5$

Compelled arcs:

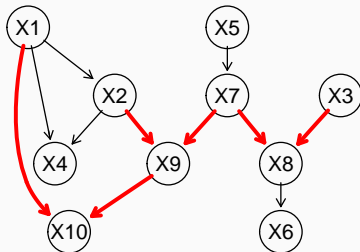
$X_1 \rightarrow X_{10}$   
 $X_2 \rightarrow X_9$   
 $X_3 \rightarrow X_8$   
 $X_7 \rightarrow X_8$   
 $X_7 \rightarrow X_9$   
 $X_8 \rightarrow X_6$   
 $X_9 \rightarrow X_{10}$

V-structures:

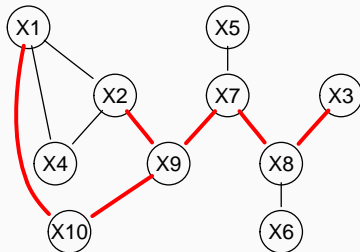
$X_1 \rightarrow X_{10} \leftarrow X_9$   
 $X_3 \rightarrow X_8 \leftarrow X_7$   
 $X_2 \rightarrow X_9 \leftarrow X_7$

# DAG, CPDAG AND EQUIVALENT DAGs

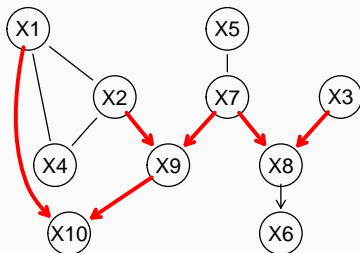
DAG



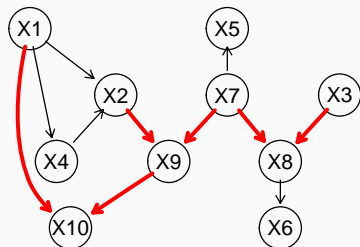
Skeleton



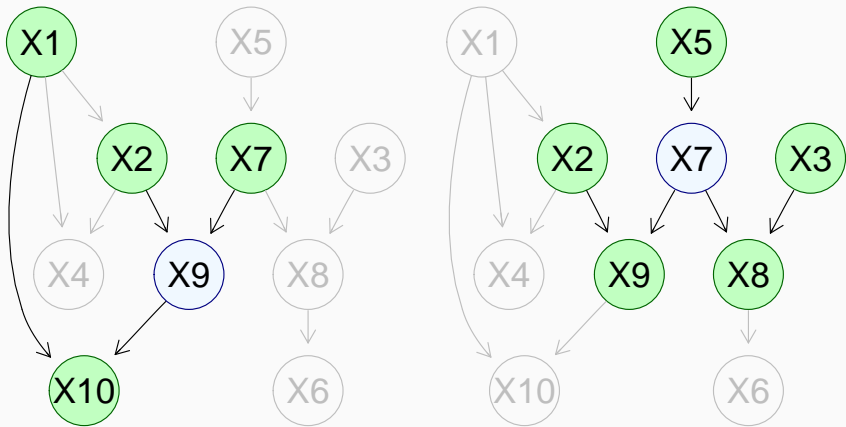
CPDAG



An Equivalent DAG



## TWO MORE EXAMPLES OF MARKOV BLANKETS





## MARKOV BLANKETS ARE SYMMETRIC

---

We can also check that **Markov blankets are symmetric**: if  $A$  is in the Markov blanket of  $B$ , then  $B$  is in the Markov blanket of  $A$ .

In which Markov blankets is  $X_9$  in?

$X_1$	$X_{10}$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$	$X_9$
TRUE	TRUE	TRUE	FALSE	FALSE	FALSE	FALSE	TRUE	FALSE	FALSE

Which nodes are in the Markov blanket of  $X_9$ ?

$X_1$	$X_{10}$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$	$X_9$
TRUE	TRUE	TRUE	FALSE	FALSE	FALSE	FALSE	TRUE	FALSE	FALSE

This is a consequence of the fact that if  $A$  is a parent of  $B$ , then  $B$  is a child of  $A$ ; and if  $A$  is a spouse of  $B$ , then  $B$  is a spouse of  $A$ .

- creating DAGs: `empty.graph()`, `set.arc()`, `drop.arc()`, `reverse.arc()`.
- model string representations: `modelstring()`, `model2network()`.
- nodes in a DAG: `nodes()`, `parents()`, `children()`, `spouses()`, `nbr()`, `mb()`.
- arcs in a DAG: `arcs()`, `path.exists()`, `dsep()`, `directed.arcs()`, `undirected.arcs()`, `compelled.arcs()`.
- DAG transformation: `subgraph()`, `moral()`, `cpdag()`.
- plotting: `graphviz.plot()`, `graphviz.compare()`.

- BNs are one of the oldest instances of machine learning models.
- BNs are a probabilistic model that use DAGs to make computations systematic in a rigorous way.
- BNs allow computer systems to perform automatically all the computations we did by hand at the beginning of this lecture.
- At the same time, BNs using DAGs means that they provide a qualitative, intuitive way to reason about complex phenomena.

Next:

- What probability distributions do we use to construct a BN?

Thanks!

Any questions?