

Explaining Bayesian networks using MAP-independence: Some new properties

Enrique Valero-Leal
Pedro Larrañaga
Concha Bielza

New Bridges between
Mathematics and Data Science

8th-11th November 2021



- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence
- 4 MAP-independence properties
- 5 MAP-independence in continuous networks
- 6 Conclusion

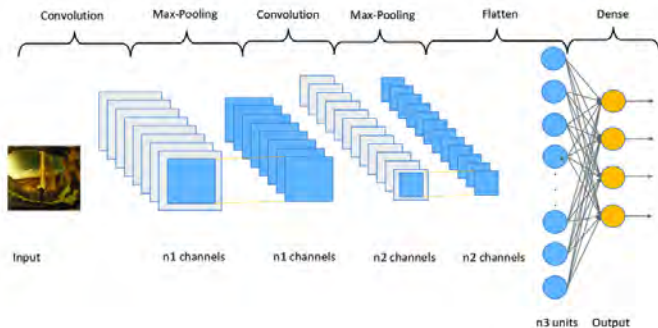
Index

- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence
- 4 MAP-independence properties
- 5 MAP-independence in continuous networks
- 6 Conclusion

Motivation for Explainable AI

Around 2012, the deep learning revolution took place

- A significant rise in complex and powerful methods
- A further insertion of such methods in our societies



Benítez-Andrades, via ResearchGate

Motivation for Explainable AI

The “reasoning” of these methods is impossible to understand

- Legal, ethical and performance problems



Explainable AI deals with better interpreting AI systems

- Miller (2019): Interpretability is the degree to which a human can understand the cause of a decision

Regarding simple models are really as simple as we believe?

- Even in interpretable models, explanations might be desirable



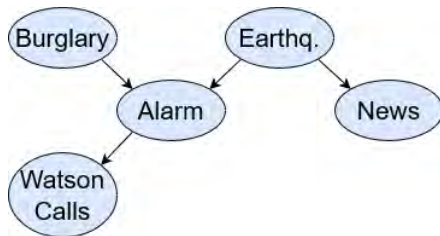
Motivation for Bayesian networks

A model based on actual mathematical and statistical theory is desirable

- Uncertainty + many variables \rightarrow intractable joint probability distribution (JPD) and inference

Bayesian networks are probabilistic graphical models that offer:

1. Modularity representing the JPD, thus compacting it
2. Easier and semantically comprehensible inference
3. Interpretable and visual relations



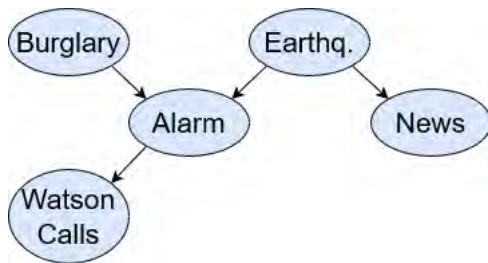
Index

- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence
- 4 MAP-independence properties
- 5 MAP-independence in continuous networks
- 6 Conclusion

Bayesian networks (Pearl, 1988) (Koller and Friedman, 2009)

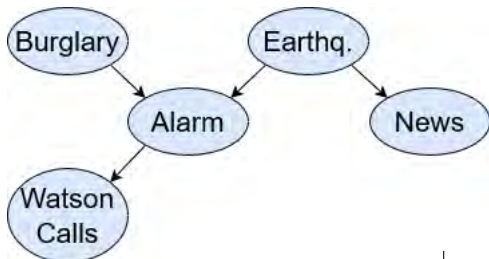
A Bayesian network $\mathcal{B} = (\mathcal{G}, \theta)$ is a probabilistic graphical model that encodes a JPD $P(X_1, \dots, X_n)$ over a set of variables $\mathbf{X} = \{X_1, \dots, X_n\}$.

Qualitatively, a Bayesian network is a directed acyclic graph $\mathcal{G} = (\mathbf{X}, A)$ that represents conditional (in)dependencies.



Quantitatively, Bayesian networks factorise the JPD $P(X_1, \dots, X_n)$ storing only the conditional probability distributions (CPD) $P(X_i | Pa_{X_i})$ of each node.

$$P(B, E, A, N, W) = P(B)P(E)P(A|B, E)P(N|E)P(W|A)$$



$P(B)$
 $P(E)$
 $P(A|B, E)$
 $P(N|E)$
 $P(W|A)$

Earthq.	News	
	t	f
f	0.1	0.9
t	0.6	0.4

Burg.	Earthq.	Alarm	
		t	f
f	f	0.05	0.95
f	t	0.4	0.6
t	f	0.9	0.1
t	t	0.99	0.01

Earthquake	t	f
	0.15	0.85



Conditional independence and d-separation

\mathbf{X}_1 and \mathbf{X}_2 are conditionally independent given \mathbf{X}_3 , $I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3)$,
iff $P(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3) = P(\mathbf{X}_1 | \mathbf{X}_3)P(\mathbf{X}_2 | \mathbf{X}_3)$.

Many independences can be identified through **d-separation**

\mathbf{X}_1 and \mathbf{X}_2 are d-separated by \mathbf{X}_3 (denoted $\mathbf{X}_1 \perp \mathbf{X}_2 | \mathbf{X}_3$) if in all undirected paths between \mathbf{X}_1 and \mathbf{X}_2 there is a node C such that:

1. C and its descendants $\notin \mathbf{X}_3$ and C is a converging connection, or
2. $C \in \mathbf{X}_3$ and C is not a converging connection.



Conditional independence and d-separation

\mathbf{X}_1 and \mathbf{X}_2 are conditionally independent given \mathbf{X}_3 , $I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3)$,
iff $P(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3) = P(\mathbf{X}_1 | \mathbf{X}_3)P(\mathbf{X}_2 | \mathbf{X}_3)$.

Many independences can be identified through **d-separation**

\mathbf{X}_1 and \mathbf{X}_2 are d-separated by \mathbf{X}_3 (denoted $\mathbf{X}_1 \perp \mathbf{X}_2 | \mathbf{X}_3$) if in all undirected paths between \mathbf{X}_1 and \mathbf{X}_2 there is a node C such that:

1. C and its descendants $\notin \mathbf{X}_3$ and C is a converging connection, or
2. $C \in \mathbf{X}_3$ and C is not a converging connection.

D-separation theorem (Verma and Pearl, 1990)

$$\mathbf{X}_1 \perp \mathbf{X}_2 | \mathbf{X}_3 \Rightarrow I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3) \quad \forall \mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3 \subseteq \mathbf{X} \text{ (disjoint)}$$

D-separation theorem (Verma and Pearl, 1990)

$$\mathbf{X}_1 \perp \mathbf{X}_2 | \mathbf{X}_3 \Rightarrow I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3) \quad \forall \mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3 \subseteq \mathbf{X} \text{ (disjoint)}$$

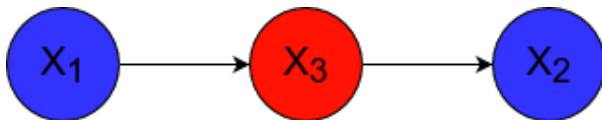
Example:



D-separation theorem (Verma and Pearl, 1990)

$$\mathbf{X}_1 \perp \mathbf{X}_2 | \mathbf{X}_3 \Rightarrow I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3) \quad \forall \mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3 \subseteq \mathbf{X} \text{ (disjoint)}$$

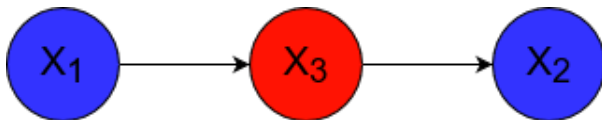
Example:



D-separation theorem (Verma and Pearl, 1990)

$$\mathbf{X}_1 \perp \mathbf{X}_2 | \mathbf{X}_3 \Rightarrow I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3) \quad \forall \mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3 \subseteq \mathbf{X} \text{ (disjoint)}$$

Example:

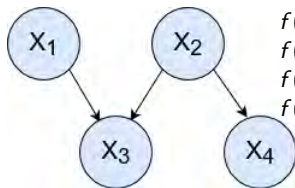


Observation: If $I(\mathbf{X}_1, \mathbf{X}_2 | \mathbf{X}_3)$ and \mathbf{X}_3 is not observed, a change in \mathbf{X}_1 is propagated to \mathbf{X}_2 via \mathbf{X}_3

Continuous Bayesian networks (Shachter and Kenley, 1989) (Geiger and Heckerman, 1994)

In continuous Bayesian networks, all the variables are continuous. One of the most widely used are linear Gaussian Bayesian networks, which factorise a multivariate Gaussian $\mathcal{N}(\boldsymbol{\mu}; \boldsymbol{\Sigma})$ into a set of linear conditional Gaussian variables.

Conditional density of Y : $f(Y|Pa_Y) = \mathcal{N}(\beta_Y + \sum_{X_i \in Pa_Y} \beta_{YX_i}(x_i - \mu_i); v_Y)$



$$f(X_1) = \mathcal{N}(\beta_1; v_1)$$

$$f(X_2) = \mathcal{N}(\beta_2; v_2)$$

$$f(X_3|X_1, X_2) = \mathcal{N}(\beta_3 + \beta_{31}(x_1 - \mu_1) + \beta_{32}(x_2 - \mu_2); v_3)$$

$$f(X_4|X_2) = \mathcal{N}(\beta_4 + \beta_{42}(x_2 - \mu_2); v_4)$$

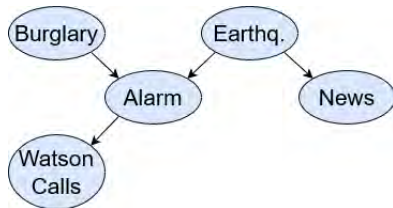
Inference

Since in a Bayesian network the JPD is factorised:

- $P(X_i) = \sum_{\mathbf{x} \setminus x_i} \prod_{j=1}^n P(X_j | Pa_{X_j})$

Probability of the alarm, $P(A)$, in the “Watson calls” network

$$P(A) = \sum_{B,E,N,W} P(B)P(E)P(A|B,E)P(N|E)P(W|A)$$



Inference: Abduction

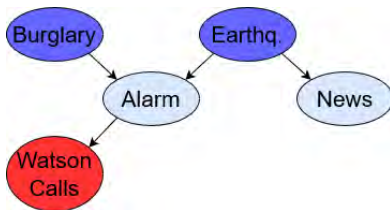
Given a set of evidence nodes $\mathbf{E} \subset \mathbf{X}$ with a value assignment \mathbf{e} and a set of explanation nodes $\mathbf{H} \subseteq \mathbf{X} \setminus \mathbf{E}$, get the most probable explanation \mathbf{h}^* for the evidence \mathbf{e} .

Maximum a posteriori (MAP)

$$\mathbf{h}^* \leftarrow \arg \max_{\mathbf{H}} P(\mathbf{H} | \mathbf{e}), \quad \mathbf{H} \subseteq \mathbf{X} \setminus \mathbf{E}$$

Example: Watson called, what is the most probable explanation of B , E ?

Answer: $\arg \max_{B,E} P(B, E | w)$



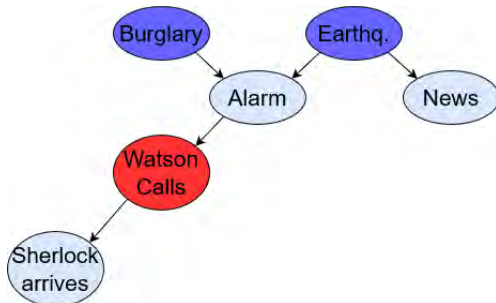
Index

- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence**
- 4 MAP-independence properties
- 5 MAP-independence in continuous networks
- 6 Conclusion

MAP-independence: Motivation

Traditionally, the variables that are conditionally dependent from the explanation set H given the evidence e are considered to be relevant for the explanation (Pearl and Paz, 1985).

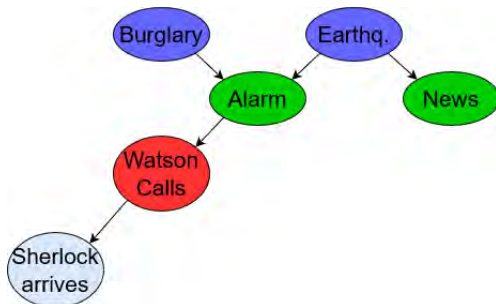
$$I(H, R|E) \Leftrightarrow R \text{ not relevant}$$



MAP-independence: Motivation

Traditionally, the variables that are conditionally dependent from the explanation set H given the evidence e are considered to be relevant for the explanation (Pearl and Paz, 1985).

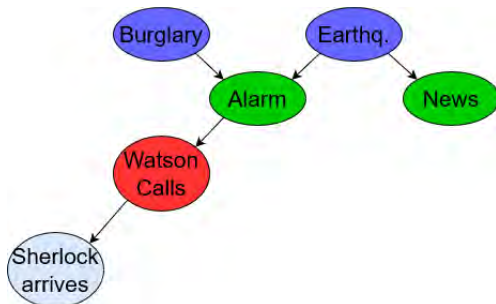
$$I(H, R|E) \Leftrightarrow R \text{ not relevant}$$



MAP-independence: Motivation

Traditionally, the variables that are conditionally dependent from the explanation set H given the evidence e are considered to be relevant for the explanation (Pearl and Paz, 1985).

$$I(H, R|E) \Leftrightarrow R \text{ not relevant}$$



Are really all conditionally dependent variables relevant in a specific context?



MAP-independence

Given a maximum a posteriori (MAP) $\mathbf{h}^* = \arg \max_{\mathbf{H}} P(\mathbf{H}|\mathbf{e})$ and a non-empty set of nodes \mathbf{R} ,

MAP-independence (Kwisthout, 2021)

Is $\forall \mathbf{r} \in \Omega(\mathbf{R}), \arg \max_{\mathbf{H}} P(\mathbf{H}, \mathbf{r}|\mathbf{e}) = \mathbf{h}^*$?

If the answer to the decision problem is Yes, then the explanation \mathbf{h}^* is MAP-independent from the set \mathbf{R} .

MAP-independence

Given a maximum a posteriori (MAP) $\mathbf{h}^* = \arg \max_{\mathbf{H}} P(\mathbf{H}|\mathbf{e})$ and a non-empty set of nodes \mathbf{R} ,

MAP-independence (Kwisthout, 2021)

Is $\forall \mathbf{r} \in \Omega(\mathbf{R}), \arg \max_{\mathbf{H}} P(\mathbf{H}, \mathbf{r}|\mathbf{e}) = \mathbf{h}^*$?

If the answer to the decision problem is Yes, then the explanation \mathbf{h}^* is MAP-independent from the set \mathbf{R} .

Observations that motivate this work:

- The properties of MAP-independence are yet to be explored.
- What about continuous variables?

Index

- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence
- 4 MAP-independence properties**
- 5 MAP-independence in continuous networks
- 6 Conclusion

MAP-independence and subsets

MAP-independence and subsets

If $R_j \subseteq R$ and h^* is MAP-independent from R , then h^* is also MAP-independent from R_j

MAP-independence and subsets

MAP-independence and subsets

If $\mathbf{R}_i \subseteq \mathbf{R}$ and \mathbf{h}^* is MAP-independent from \mathbf{R} , then \mathbf{h}^* is also MAP-independent from \mathbf{R}_i

Let $\mathbf{R}_j = \mathbf{R} \setminus \mathbf{R}_i$:

$\forall (\mathbf{r}_i, \mathbf{r}_j) \in \Omega(\mathbf{R}_i) \times \Omega(\mathbf{R}_j)$, $\arg \max_{\mathbf{H}} P(\mathbf{H}, \mathbf{r}_i, \mathbf{r}_j | \mathbf{e}) = \mathbf{h}^*$

Make the *arg max* explicit.

$\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $\forall (\mathbf{r}_i, \mathbf{r}_j) \in \Omega(\mathbf{R}_i) \times \Omega(\mathbf{R}_j)$, $P(\mathbf{h}^*, \mathbf{r}_i, \mathbf{r}_j | \mathbf{e}) > P(\bar{\mathbf{h}}, \mathbf{r}_i, \mathbf{r}_j | \mathbf{e})$

MAP-independence and subsets

MAP-independence and subsets

If $R_i \subseteq R$ and h^* is MAP-independent from R , then h^* is also MAP-independent from R_i

Let $R_j = R \setminus R_i$:

$$\forall (r_i, r_j) \in \Omega(R_i) \times \Omega(R_j), \arg \max_{\mathbf{H}} P(\mathbf{H}, r_i, r_j | \mathbf{e}) = h^*$$

Make the *arg max* explicit.

$$\forall \bar{h} \in \Omega(\mathbf{H}) \setminus h^*, \forall (r_i, r_j) \in \Omega(R_i) \times \Omega(R_j), P(h^*, r_i, r_j | \mathbf{e}) > P(\bar{h}, r_i, r_j | \mathbf{e})$$

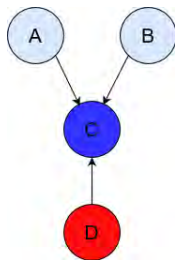
Marginalize $P(h^*, r_i | \mathbf{e})$

$$\forall \bar{h} \in \Omega(\mathbf{H}) \setminus h^*, \forall r_i \in \Omega(R_i), \sum^{R_j} P(h^*, r_i, r_j | \mathbf{e}) > \sum^{R_j} P(\bar{h}, r_i, r_j | \mathbf{e})$$

$$\forall \bar{h} \in \Omega(\mathbf{H}) \setminus h^*, \forall r_i \in \Omega(R_i), P(h^*, r_i | \mathbf{e}) > P(\bar{h}, r_i | \mathbf{e})$$

Introduce *arg max*

$$\forall r_i \in \Omega(R_i), \arg \max_{\mathbf{H}} P(\mathbf{H}, r_i | \mathbf{e}) = h^*$$

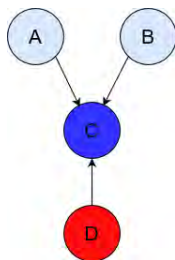


		$P(A)$
		A
t		f
ε		$1 - \epsilon$

		$P(B)$
		B
t		f
μ		$1 - \mu$

		$P(C A, B, d)$	
		C	
A	B	t	f
f	f	0.55	0.45
f	t	0.6	0.4
t	f	0.65	0.35
t	t	0.99	0.01

MAP query: If d , most probable value assignment for C ?
 $\arg \max_C P(C|d) = c$



$P(A)$	
	A
t	f
ε	$1 - \epsilon$

$P(B)$	
	B
t	f
μ	$1 - \mu$

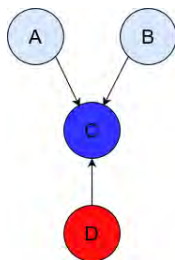
		$P(C A, B, d)$	
A	B	C	
		t	f
f	f	0.55	0.45
f	t	0.6	0.4
t	f	0.65	0.35
t	t	0.99	0.01

MAP query: If d , most probable value assignment for C ?

$$\arg \max_C P(C|d) = c$$

Check MAP-independence of the explanation c from A and B

$$\forall (a, b) \in \Omega(A) \times \Omega(B), \arg \max_C P(C|A, B, d) = \arg \max_C P(C|d) = c$$



$P(A)$	
	A
t	f
ε	1 - ε

$P(B)$	
	B
t	f
μ	1 - μ

		$P(C A, B, d)$	
		C	
A	B	t	f
f	f	0.55	0.45
f	t	0.6	0.4
t	f	0.65	0.35
t	t	0.99	0.01

MAP query: If d , most probable value assignment for C ?

$$\arg \max_C P(C|d) = c$$

Check MAP-independence of the explanation c from A and B

$$\forall (a, b) \in \Omega(A) \times \Omega(B), \arg \max_C P(C|A, B, d) = \arg \max_C P(C|d) = c$$

Check MAP-independence of c from just A

$$\forall a \in \Omega(A), \arg \max_C P(C|a, d)$$

$$P(C|a, d) = \mu \cdot 0.99 + (1 - \mu) \cdot 0.65 > 0.5; \arg \max_C P(C|a, d) = c$$

$$P(C|\bar{a}, d) = \mu \cdot 0.6 + (1 - \mu) \cdot 0.55 > 0.5; \arg \max_C P(C|\bar{a}, d) = c$$

MAP-independence and conditional independence

MAP-independence and conditional independence

If $I(\mathbf{H}, \mathbf{R}_j | \mathbf{R}_i)$ and \mathbf{h}^* is MAP-independent from \mathbf{R}_i , then \mathbf{h}^* is MAP-independent from \mathbf{R}_j .

MAP-independence and conditional independence

MAP-independence and conditional independence

If $I(\mathbf{H}, \mathbf{R}_j | \mathbf{R}_i)$ and \mathbf{h}^* is MAP-independent from \mathbf{R}_i , then \mathbf{h}^* is MAP-independent from \mathbf{R}_j .

Since $P(\mathbf{h}, \mathbf{r}_j | \mathbf{e}) \propto P(\mathbf{h} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_j \in \Omega(\mathbf{R}_j)$, $\mathbf{h}^* = \arg \max_{\mathbf{H}} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_j \in \Omega(\mathbf{R}_j)$, $\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) > P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e})$

MAP-independence and conditional independence

MAP-independence and conditional independence

If $I(\mathbf{H}, \mathbf{R}_j | \mathbf{R}_i)$ and \mathbf{h}^* is MAP-independent from \mathbf{R}_i , then \mathbf{h}^* is MAP-independent from \mathbf{R}_j .

Since $P(\mathbf{h}, \mathbf{r}_i | \mathbf{e}) \propto P(\mathbf{h} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_i \in \Omega(\mathbf{R}_i)$, $\mathbf{h}^* = \arg \max_{\mathbf{H}} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_i \in \Omega(\mathbf{R}_i)$, $\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) > P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e})$

For every probability distribution $P(\mathbf{R}_i)$, the following holds

$\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $\sum^{R_i} P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i) > \sum^{R_i} P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i)$

$\arg \max_{\mathbf{H}} \sum^{R_i} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i) = \mathbf{h}^*$

MAP-independence and conditional independence

MAP-independence and conditional independence

If $I(\mathbf{H}, \mathbf{R}_j | \mathbf{R}_i)$ and \mathbf{h}^* is MAP-independent from \mathbf{R}_i , then \mathbf{h}^* is MAP-independent from \mathbf{R}_j .

Since $P(\mathbf{h}, \mathbf{r}_i | \mathbf{e}) \propto P(\mathbf{h} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_i \in \Omega(\mathbf{R}_i)$, $\mathbf{h}^* = \arg \max_{\mathbf{H}} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_i \in \Omega(\mathbf{R}_i)$, $\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) > P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e})$

For every probability distribution $P(\mathbf{R}_i)$, the following holds

$\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $\sum^{R_i} P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i) > \sum^{R_i} P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i)$

$\arg \max_{\mathbf{H}} \sum^{R_i} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i) = \mathbf{h}^*$

i.e. any probability distribution $P(\mathbf{R}_i)$ over \mathbf{R}_i cannot alter \mathbf{h}^*

MAP-independence and conditional independence

MAP-independence and conditional independence

If $I(\mathbf{H}, \mathbf{R}_j | \mathbf{R}_i)$ and \mathbf{h}^* is MAP-independent from \mathbf{R}_i , then \mathbf{h}^* is MAP-independent from \mathbf{R}_j .

Since $P(\mathbf{h}, \mathbf{r}_i | \mathbf{e}) \propto P(\mathbf{h} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_i \in \Omega(\mathbf{R}_i)$, $\mathbf{h}^* = \arg \max_{\mathbf{H}} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e})$

$\forall \mathbf{r}_i \in \Omega(\mathbf{R}_i)$, $\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) > P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e})$

For every probability distribution $P(\mathbf{R}_i)$, the following holds

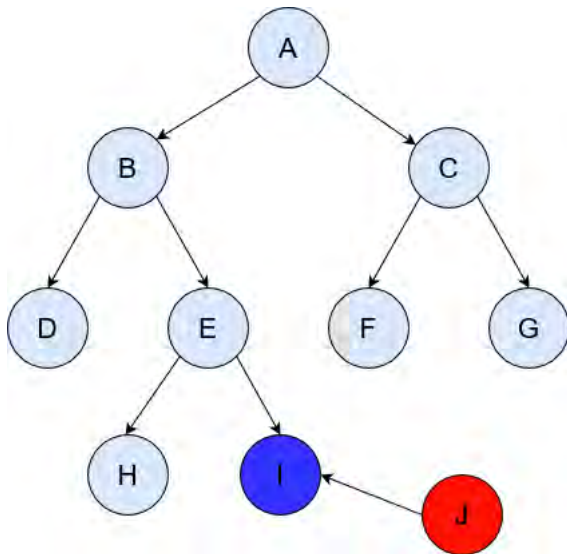
$\forall \bar{\mathbf{h}} \in \Omega(\mathbf{H}) \setminus \mathbf{h}^*$, $\sum^{R_i} P(\mathbf{h}^* | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i) > \sum^{R_i} P(\bar{\mathbf{h}} | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i)$

$\arg \max_{\mathbf{H}} \sum^{R_i} P(\mathbf{H} | \mathbf{r}_i, \mathbf{e}) P(\mathbf{r}_i) = \mathbf{h}^*$

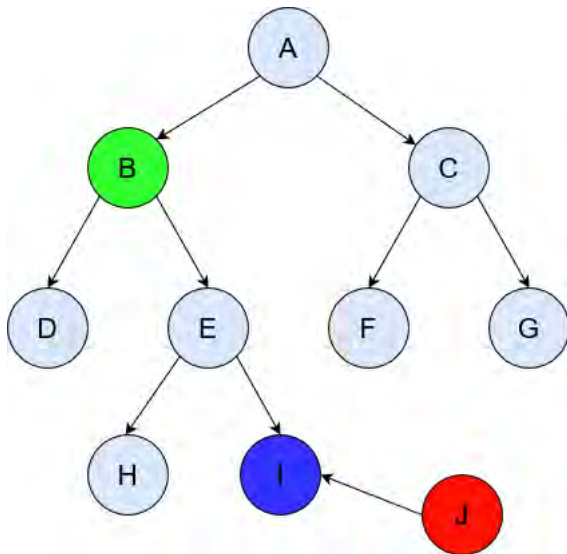
i.e. any probability distribution $P(\mathbf{R}_i)$ over \mathbf{R}_i cannot alter \mathbf{h}^*

Since \mathbf{R}_j only affects \mathbf{H} through \mathbf{R}_i (c.i. given \mathbf{R}_i), \mathbf{R}_j cannot alter \mathbf{h}^*

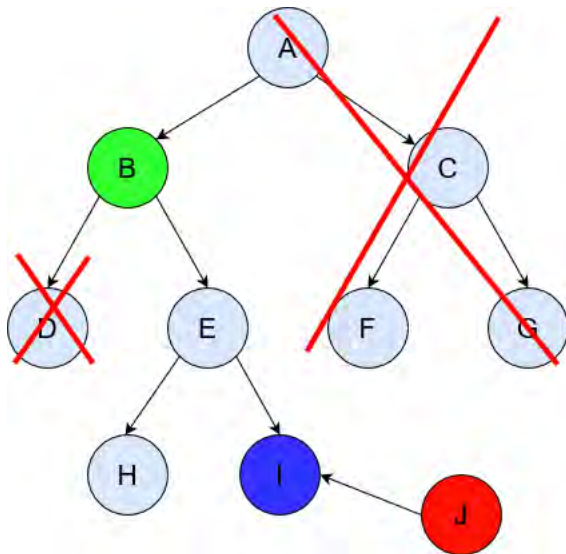
Evidence **J**, MAP query over **I** yields i^*



i^* is MAP-independent from **B**



i^* is also MAP-independent from the nodes d-separated from I given B



Index

- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence
- 4 MAP-independence properties
- 5 **MAP-independence in continuous networks**
- 6 Conclusion

MAP-independence in continuous Bayesian networks

MAP-independence

Is $\forall \mathbf{r} \in \Omega(\mathbf{R}) \arg \max_{\mathbf{H}} P(\mathbf{H}, \mathbf{r} | \mathbf{e}) = \mathbf{h}^*$?

arg max of a continuous distribution is a real number or a vector of real numbers (the mode or a vectorial mode).

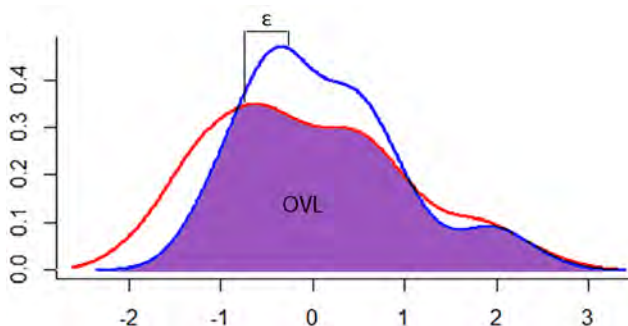
- Use a threshold of tolerance ϵ to compare the modes
- Another option is to compare the probability distributions directly.
- Both ideas can be captured within an abstract function *dist()*

Change in the formulation

- Is $\forall \mathbf{r} \in \Omega(\mathbf{R}) \text{dist}(P(\mathbf{H} | \mathbf{e}), P(\mathbf{H} | \mathbf{e}, \mathbf{r})) < \epsilon$?

Some examples of $dist()$ functions (we could define others):

1. $\sum(|\arg \max_{\mathbf{H}} P(\mathbf{H}|\mathbf{e}) - \arg \max_{\mathbf{H}} P(\mathbf{H}|\mathbf{e}, \mathbf{r})|)$.
2. The Jensen-Shannon divergence, $JSD(P(\mathbf{H}|\mathbf{e})||P(\mathbf{H}|\mathbf{e}, \mathbf{r}))$.
3. The overlapping coefficient between the two distributions, $OVL(P(\mathbf{H}|\mathbf{e}), P(\mathbf{H}|\mathbf{e}, \mathbf{r}))$.



ϵ -MAP-independence

MAP-independence

Is $\forall \mathbf{r} \in \Omega(\mathbf{R}) \arg \max_{\mathbf{H}} P(\mathbf{H}, \mathbf{r} | \mathbf{e}) = \mathbf{h}^*$?

Iterating over all values of \mathbf{R} ($\forall \mathbf{r} \in \Omega(\mathbf{R})$) \rightarrow not possible (continuous).

- Sample from $P(\mathbf{R})$ and check if on average the condition holds.

Merging the two proposals:

ϵ -MAP-independence

Is $\text{dist}(P(\mathbf{H} | \mathbf{e}), P(\mathbf{H} | \mathbf{e}, \mathbf{r})) < \epsilon$, where \mathbf{r} is sampled from $P(\mathbf{R})$?

Gaussian linear Bayesian networks

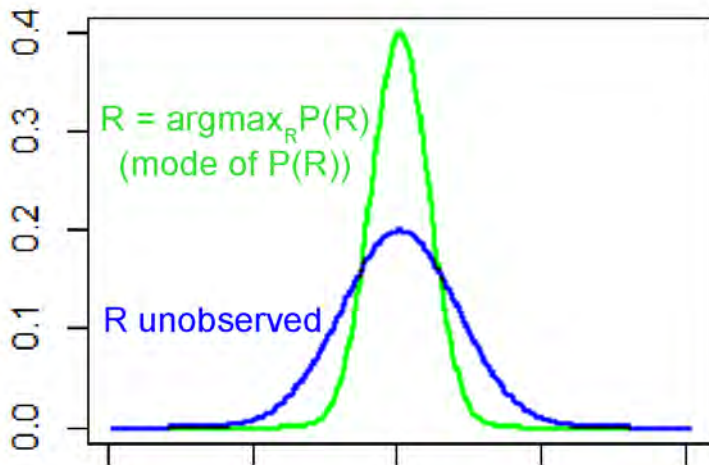
In the particular case of Gaussian Bayesian networks, we can take advantage of the linear nature of the node means.

Conditional density of a node Y with parents Pa_Y :

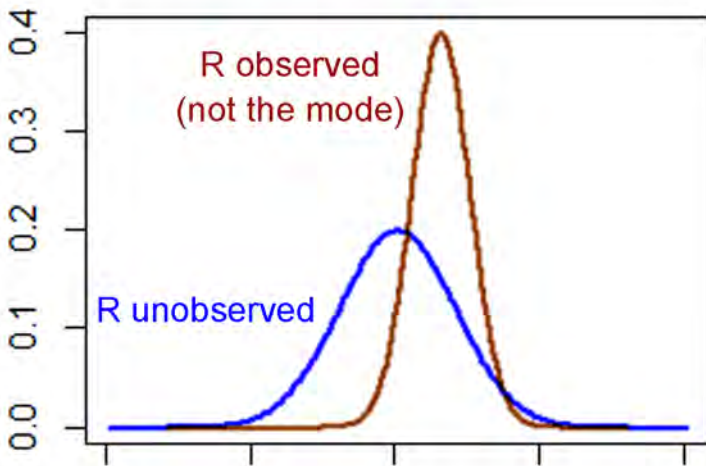
$$f(Y|Pa_Y) = \mathcal{N}(\beta_Y + \sum_{X_i \in Pa_Y} \beta_{YX_i}(x_i - \mu_i) ; v_Y)$$

Some prior observation regarding inference:

If $r = \arg \max_R P(\mathbf{R})$, then $\text{dist}(P(\mathbf{H}|\mathbf{e}), P(\mathbf{H}|\mathbf{e}, r))$ is minimum.

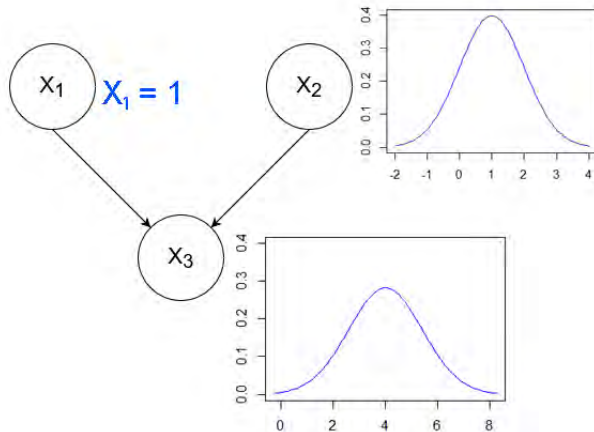


If $r = \arg \max_{\mathbf{R}} P(\mathbf{R})$, then $\text{dist}(P(\mathbf{H}|\mathbf{e}), P(\mathbf{H}|\mathbf{e}, r))$ is minimum.



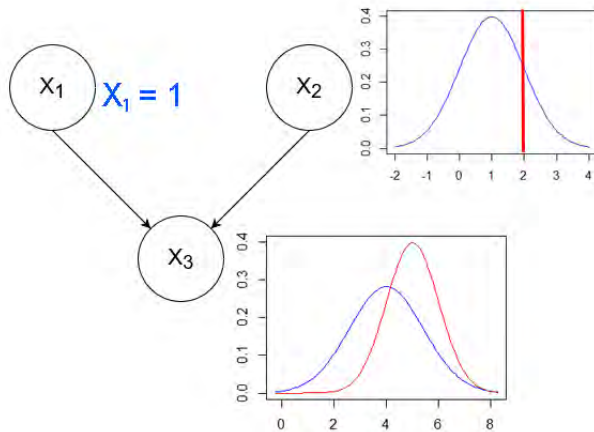
Gaussian linear Bayesian networks

We observe a certain evidence $X_1 = 1$. Observe $P(X_2)$ and $P(X_3|X_1 = 1)$



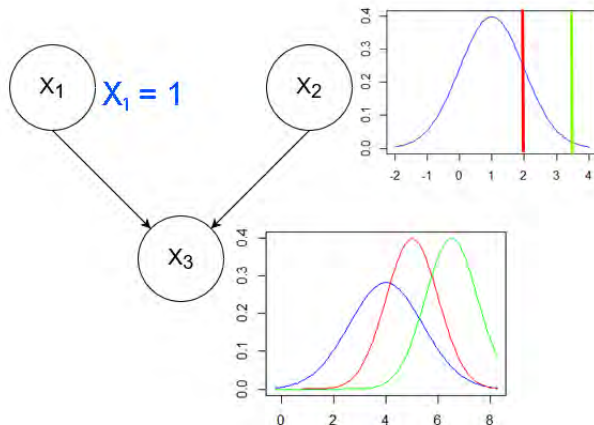
Gaussian linear Bayesian networks

We sample "2" from $P(X_2)$. In red, $P(X_3|X_1 = 1, X_2 = 2)$



Gaussian linear Bayesian networks

Since $X_2 = 3.5$ is further from the X_2 mode ($\arg \max_{X_2} P(X_2) = 1$) than $X_2 = 2$, then $P(X_3|X_1 = 1, X_2 = 3.5)$ will be more dissimilar (in green)



ϵ -MAP-independence

Is $\text{dist}(P(\mathbf{H}|\mathbf{e}), P(\mathbf{H}|\mathbf{e}, \mathbf{r})) < \epsilon$, where \mathbf{r} is sampled from $P(\mathbf{R})$?

A single sample \mathbf{r}' from $P(\mathbf{R})$ gives us plenty of information

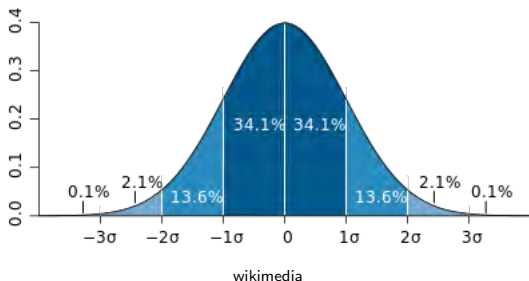
- If \mathbf{r}' meets the ϵ -MAP-independence criterion, any sample in range $[\arg \max_{\mathbf{R}} P(\mathbf{R}), \mathbf{r}']$ will.
- Otherwise, not all values in range $[\arg \max_{\mathbf{R}} P(\mathbf{R}), \mathbf{r}']$ meet it.

The question remains:

- How many and which \mathbf{r}' 's to select?

We want a value that gives us information about a large part of $P(\mathbf{R})$

- Checking only 2σ and -2σ let us check the 95% of its probability.
- We need $2^{|\mathbf{R}|}$ samples (combinations of $\pm 2\sigma$ for each $R \in \mathbf{R}$).



Index

- 1 Introduction
- 2 Background
 - Bayesian networks
 - Conditional independence and d-separation
 - Inference
- 3 MAP-independence
- 4 MAP-independence properties
- 5 MAP-independence in continuous networks
- 6 Conclusion



Conclusion

The interest on MAP-independence is two-fold:

1. It lets us study the node relevance and explanation robustness with a stricter criterion.
2. The methodology is extremely easy to understand → XAI principle, legislation...

Although our extension to the continuous domain relies heavily on heuristics, simulation methods are still useful and will enable MAP-independence to be applied in many more domains.

Future ways

Most of the datasets nowadays are either continuous or contain continuous and discrete variables

- Extend the methodology to hybrid Bayesian networks.

We can use the formulated MAP-independence properties to design an efficient algorithm

- Start from the closest nodes to H and explore “outwards” (i.e., starting from the Markov blanket).
- Start from a reduced number of variables.

The definition of MAP-independence is related to counterfactuals (Molnar, 2020)

- Check if a change in R produces a counterfactual in H .

Acknowledgements

Acknowledgements:

- Spanish Ministry of Science and Innovation through the PID 2019-109247GB-I00 project,
- The BBVA Foundation (2019 Call) through the "Score-based nonstationary temporal Bayesian networks. Applications in climate and neuroscience" project,
- The BBVA Foundation's grants (2020 Call) for Scientific Investigation Teams SARS-CoV-2 and COVID-19 through the "Outcome prediction and treatment efficiency inpatients hospitalized with Covid-19 in Madrid: A Bayesian network approach" project.

References I

- Geiger, D. and Heckerman, D. (1994). Learning Gaussian networks. In *Proceedings of the Tenth international conference on Uncertainty in artificial intelligence*, pages 235–243. Elsevier.
- Koller, D. and Friedman, N. (2009). *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press.
- Kwisthout, J. (2021). Explainable AI Using MAP-independence. In *European Conference on Symbolic and Quantitative Approaches with Uncertainty*, pages 243–254. Springer.
- Miller, T. (2019). Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38.
- Molnar, C. (2020). *Interpretable Machine Learning*. Lulu. com.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann.
- Pearl, J. and Paz, A. (1985). *Graphoids: A graph-based logic for reasoning about relevance relations*. University of California (Los Angeles). Computer Science Department.
- Shachter, R. D. and Kenley, C. R. (1989). Gaussian influence diagrams. *Management Science*, 35(5):527–550.
- Verma, T. and Pearl, J. (1990). Causal networks: semantics and expressiveness. In *Proceedings of the Fourth Annual Conference on Uncertainty in Artificial Intelligence*, pages 69–78.