

Marta Alonso Tubía

20 Julio 2023

Tutores: Concha Bielza & Pedro Larrañaga



- 

## 1 Explicaciones en redes Bayesianas

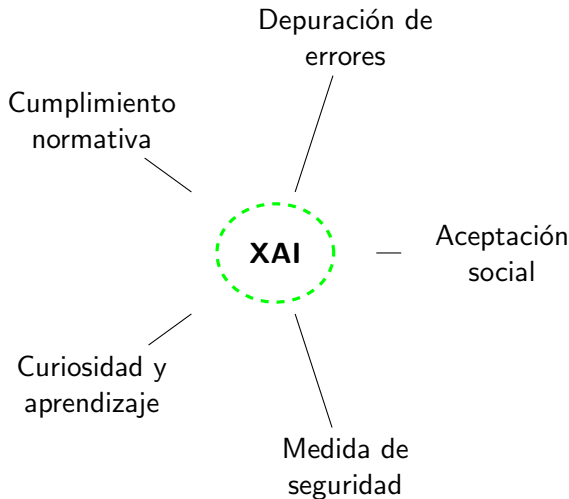
### 3 Inferencia abductiva

#### 4 Explicación más relevante

## 5 Conclusiones



# Motivación



## Definición

Una red Bayesiana (RB), es una tupla  $\mathcal{B} = (\mathcal{G}, \theta)$ , donde  $\mathcal{G} = (V, A)$  es un grafo dirigido y acíclico con un conjunto de nodos  $V = \{X_1, \dots, X_n\}$  y un conjunto de arcos  $A \subseteq V \times V$ . El conjunto  $\theta = \{P(x_i | \mathbf{X}_{Pa_{X_i}})\}$  define una CPD para cada nodo del grafo. Una RB representa una distribución conjunta de probabilidad  $P(\mathbf{x})$  del vector aleatorio  $\mathbf{X} = (X_1, \dots, X_n)$

$$P(\mathbf{x}) = \prod_{i=1}^n P(x_i | \mathbf{X}_{Pa_{X_i}}).$$

Adecuadas para interpretabilidad:

- Representa visualmente las dependencias entre variables.
- Modela la incertidumbre usando probabilidades.
- Permite aprendizaje incremental.
- Permite incorporar conocimientos previos.



Inferencia {

- Probabilidad marginal:  $P(\mathbf{Y})$
- Probabilidad condicionada:  $P(\mathbf{Y}|\mathbf{e})$
- Inferencia abductiva:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{Y}} P(\mathbf{Y}|\mathbf{e}), \mathbf{Y} \subseteq \mathbf{X}$
- Explicación más relevante (MRE)



- 

## Algoritmos exactos

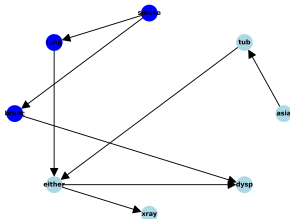
- Eliminación de variables y Árbol de cliqués.
- **Objetivo:** generar visualizaciones que ilustren pasos intermedios y justifiquen decisiones del razonamiento.
- Librerías: networkX y matplotlib.



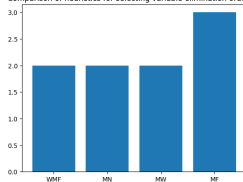


## Eliminación de variables

- Eliminación sistemática y ordenada de variables en una lista de factores.
- Es posible una optimización inicial: **pruning**.
- Ordenación clave en la eficiencia del procedimiento.
  - Comparación (heurísticas) usando *induced-width*.



Comparison of heuristics for selecting variable elimination ordering



Variable a eliminar Earthquake  
 Factores que participan

Alarm	Burglary	Earthquake	phi (Alarm,Burglary,Earthquake)
Alarm(0)	Burglary(0)	Earthquake(0)	0.9990
Alarm(0)	Burglary(0)	Earthquake(1)	0.7100
Alarm(0)	Burglary(1)	Earthquake(0)	0.0600
Alarm(0)	Burglary(1)	Earthquake(1)	0.0500
Alarm(1)	Burglary(0)	Earthquake(0)	0.0010
Alarm(1)	Burglary(0)	Earthquake(1)	0.2900
Alarm(1)	Burglary(1)	Earthquake(0)	0.9400
Alarm(1)	Burglary(1)	Earthquake(1)	0.9500
Earthquake	phi (Earthquake)		
Earthquake(0)	0.9980		
Earthquake(1)	0.0020		



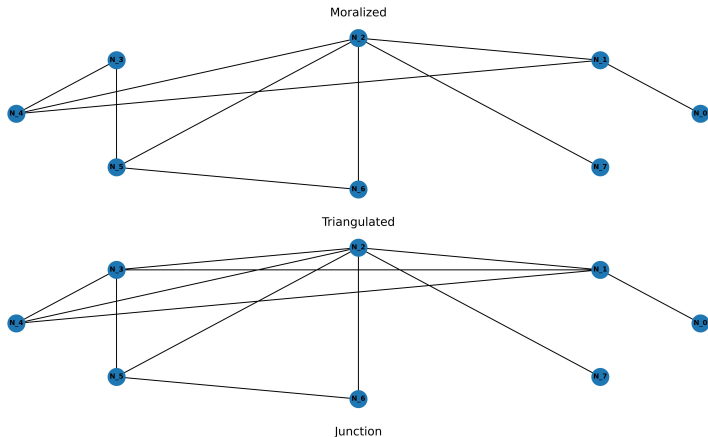
Phi

Burglary	Alarm	phi (Burglary,Alarm)
Burglary(0)	Alarm(0)	0.9984
Burglary(0)	Alarm(1)	0.0016
Burglary(1)	Alarm(0)	0.0600
Burglary(1)	Alarm(1)	0.9400



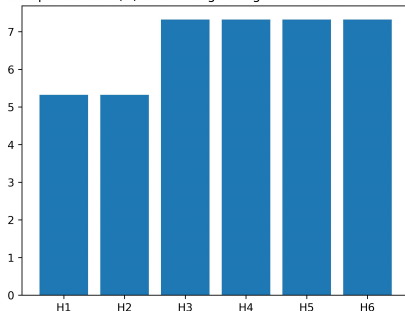
# Árbol de cliques

## Construcción del árbol de cliques

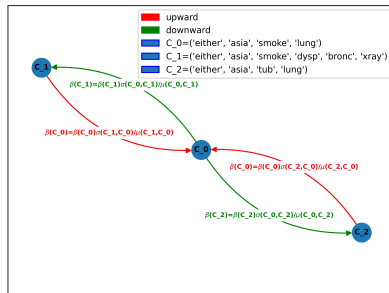


## Coste de triangulación y message passing

Comparison of  $w(G)$  after triangulating with different heuristics



$$w(\mathcal{G}) = \log_2 \sum_C \prod_{v_i \in C} n_i$$



## Algoritmos aproximados: Muestreo de Gibbs

Repetir muestreo hasta que la distribución conjunta esté más cerca de la distribución posterior  $P(\mathbf{X}|\mathbf{e})$ .

### ¿CONVERGENCIA?

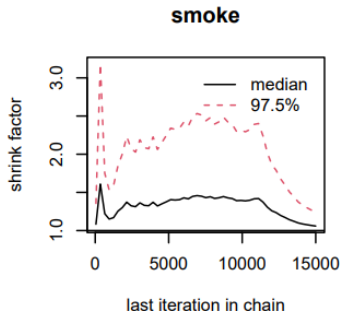
- 1 Opciones para la convergencia del muestreo de Gibbs → métodos estadísticos de diagnóstico de convergencia.
- 2 Limitaciones de librerías que hacen inferencia aproximada (*bnlearn*, *pgmpy*, *pomegranate...*).  
→ extender *pgmpy* con otra librería especializada en test de diagnóstico de convergencia de cadenas de Markov (*coda(R)*).
- 3 Comunicar python-R y transformar los objetos de una a otra librería. → *rpy2*. Implementar.



## Gelman-Rubin

	Point est
asia	1.00
tub	1.32
either	1.85
dysp	1.04
smoke	1.03
bronc	1.00
lung	1.73
xray	1.55

Multivariate psrf: 1.44



## Índice

- 1 Explicaciones en redes Bayesianas
  - Motivación
  - Tipos de inferencia
- 2 Probabilidad marginal/condicionada
- 3 Inferencia abductiva
- 4 Explicación más relevante
- 5 Conclusiones



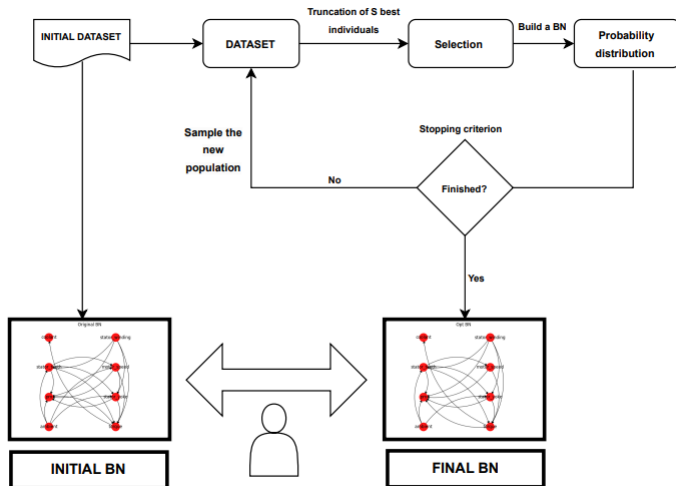
# Propuesta

- Librería: **EDAspy**. Implementación: extensión de **EGNA**.





# Metodología



- **Procedimiento:** minimizar la función de coste  $C_l$ ,  $l = 1, 2, 3, \dots$

$$C_l(\mathbf{x}) = \begin{cases} \log(f_l(\boldsymbol{\mu}_l)) - \log(f_l(\mathbf{x})) & \|\tilde{\mathbf{e}} - \mathbf{e}\|_\infty \leq \epsilon \\ 99999 & \text{otro caso} \end{cases}$$

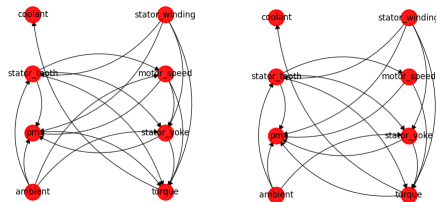
donde  $\mathbf{x} = (\mathbf{y}, \tilde{\mathbf{e}})$ , con criterio de parada:

$$\text{death iter} \geq N_{iter} \quad \text{or} \quad \frac{f_l(\mu_l)}{f_0(\mu_0)} \geq \alpha > 1$$



## Resultados (dataset motor eléctrico)

**Objetivo:** mantenimiento del motor más eficiente.



- Incremento en la función de densidad 1.638
- Las dependencias directas velocidad motor - temperatura estator y velocidad motor - temperatura ambiente, desaparecen.
- La temperatura del motor decrece su correlación notablemente con la temperatura de diferentes partes del estator.



- 

## Explicación más relevante (MRE)

- Una buena explicación es **concisa** y **precisa** → necesitamos una métrica.

$$GBF(\mathbf{x}; \mathbf{e}) \equiv \frac{P(\mathbf{e}|\mathbf{x})}{P(\mathbf{e}|\bar{\mathbf{x}})}$$

El GBF presenta características teóricas que permiten identificar automáticamente la explicación más relevante.

### Definición

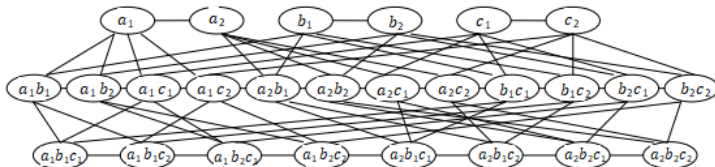
*Sea  $\mathbf{M}$  un conjunto de variables target, y  $\mathbf{e}$  la evidencia parcial en el resto de variables en una RB. La explicación más relevante es el problema de encontrar una explicación  $\mathbf{x}$  para  $\mathbf{e}$  de forma que tenga el máximo factor generalizado de Bayes  $GBF(\mathbf{x}; \mathbf{e})$ , i.e.,*

$$MRE(\mathbf{M}; \mathbf{e}) \equiv \operatorname{argmax}_{\mathbf{x}, \emptyset \subset \mathbf{x} \subseteq \mathbf{M}} GBF(\mathbf{x}; \mathbf{e}).$$



# Implementación y k-MRE

- **Implementación:** Forward search

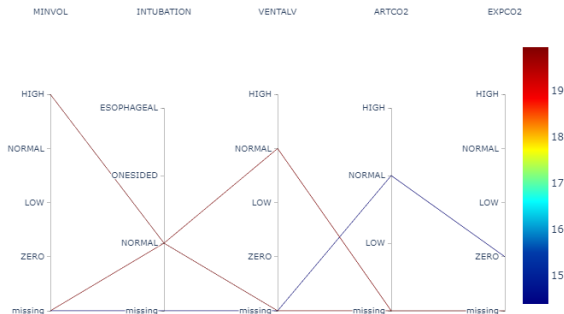


- **k-MRE:** poder escoger entre posibles explicaciones y comparar la calidad de las mismas. → Diversidad y GBF



# Resultados (dataset ALARM)

Evidencia: baja ventilación pulmonar.



MINVOL	INTUBATION	VENTLV	ARTCO2	EXPCO2	GBF	Hamming distance
HIGH	NORMAL	missing	missing	missing	19.93	0 (MRE)
missing	NORMAL	NORMAL	missing	missing	19.79	2
missing	missing	missing	NORMAL	ZERO	14.39	4



- 



- 

- 



