

# REDES BAYESIANAS

## CAP 14 (14.1 - 14.3)

---

Parcialmente adaptado de  
<http://aima.eecs.berkeley.edu>

# Resumo

- Sintaxe
- Semântica
- Distribuições parametrizadas

# Redes Bayesianas

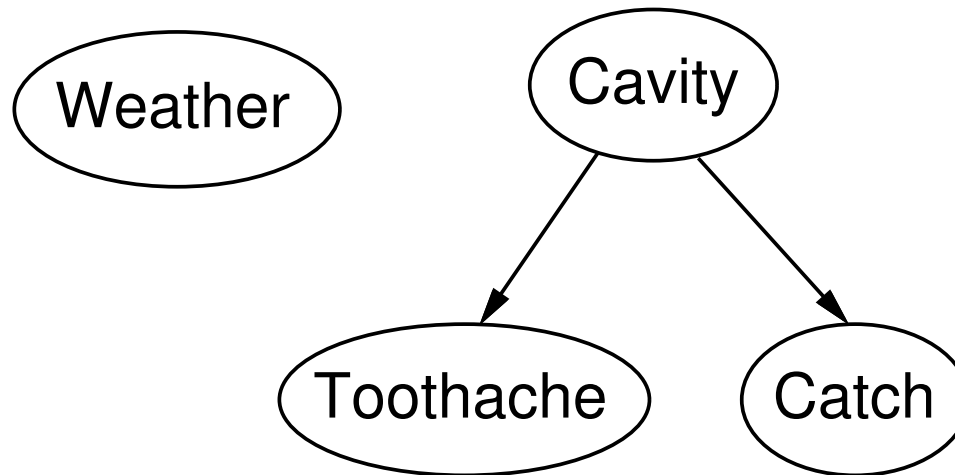
- Notação gráfica simples para **asserções de independência condicional**, e portanto para uma especificação compacta de distribuições conjuntas totais
- Sintaxe:
  - um conjunto de nós, um por variável
  - um grafo acíclico dirigido (arco  $\approx$  “influencia directamente”)
  - uma distribuição condicional para cada nó dados os seus pais:

$$\mathbf{P}(X_i \mid \text{Parents}(X_i))$$

- No caso mais simples, distribuição condicional representada como uma tabela de probabilidade condicionada (CPT) especificando a distribuição de  $X_i$  para cada combinação de valores dos pais

# Exemplo

- A topologia da rede representa asserções de independência condicional:



- *Weather* é independente das outras variáveis
- *Toothache* e *Catch* são condicionalmente independentes dado *Cavity*

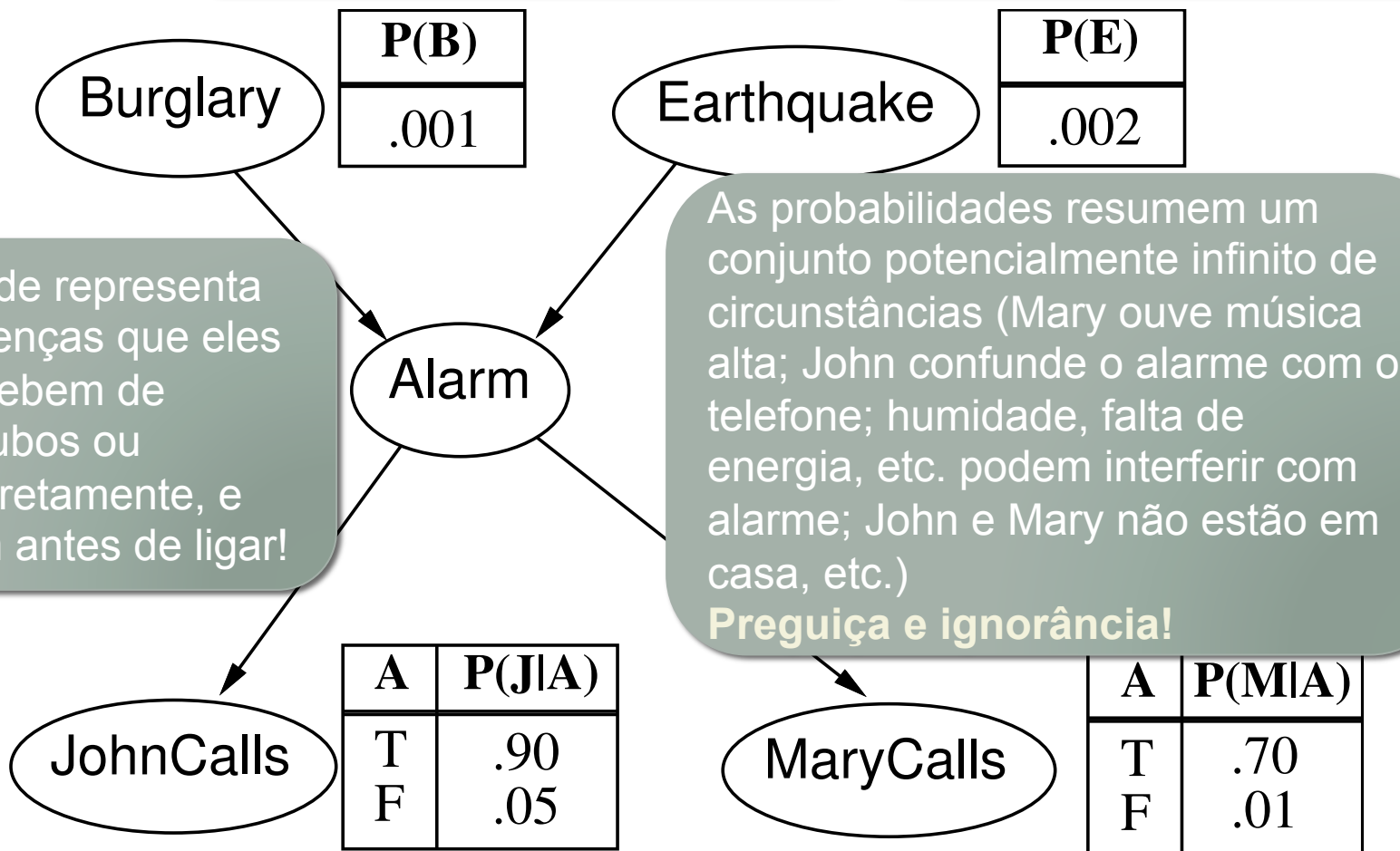
# Exemplo

- Estou no trabalho, o vizinho John telefona-me para dizer que o meu alarme está a tocar, mas a vizinha Mary não telefona. Ocasionalmente dispara por causa de pequenos tremores de terra. A minha casa esta a ser assaltada?
- Variáveis: *Burglar*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*
- Topologia da rede reflecte conhecimento “causal”:
  - Um assaltante pode fazer disparar o alarme
  - Um tremor de terra pode fazer disparar o alarme
  - O alarme pode fazer com que a Mary telefone
  - O alarme pode fazer com que o John telefone

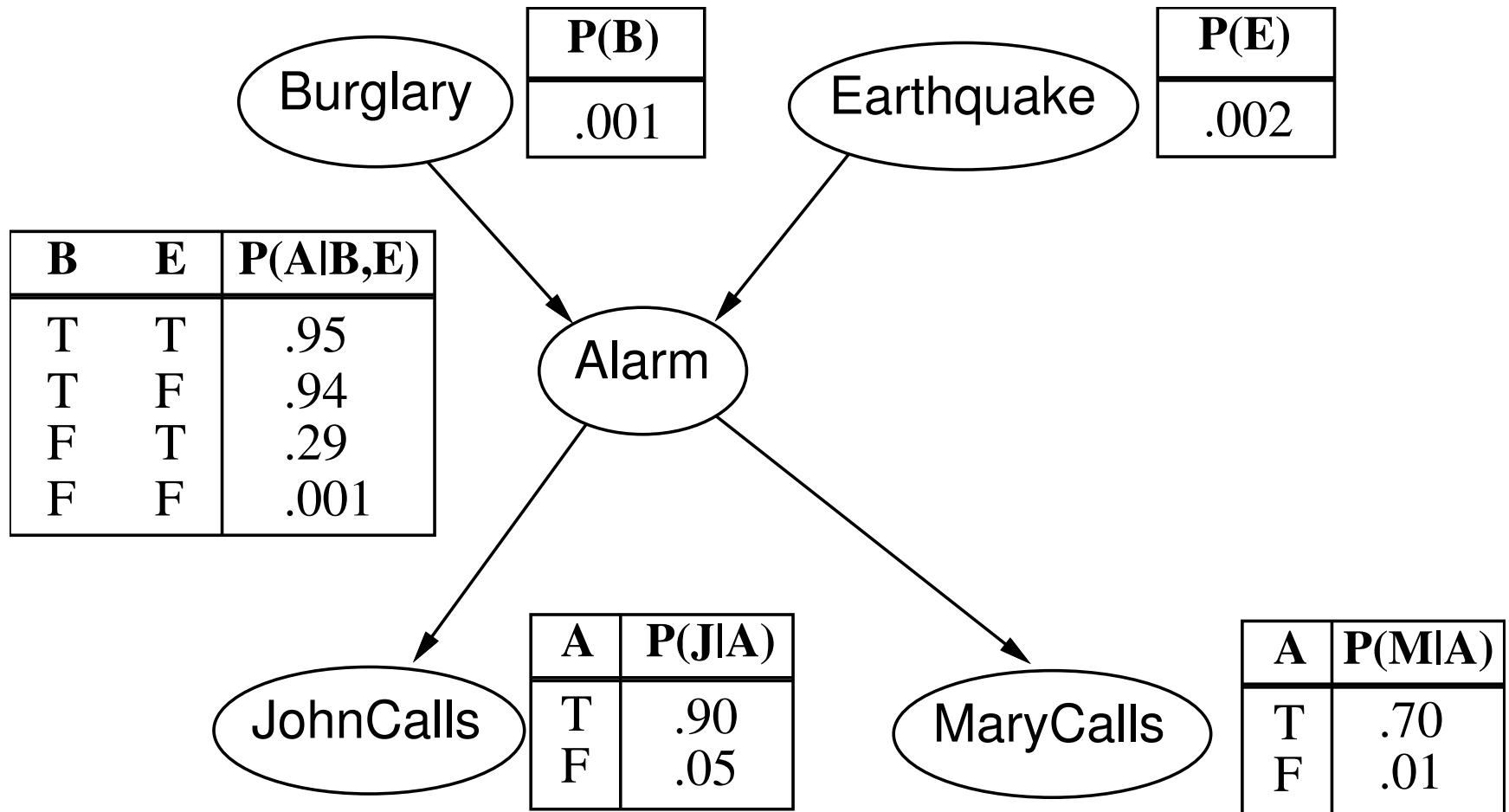
# Exemplo

Roubos e terremotos afectam diretamente a probabilidade do alarme tocar

O Facto de John e Mary telefonarem só depende diretamente do alarme.

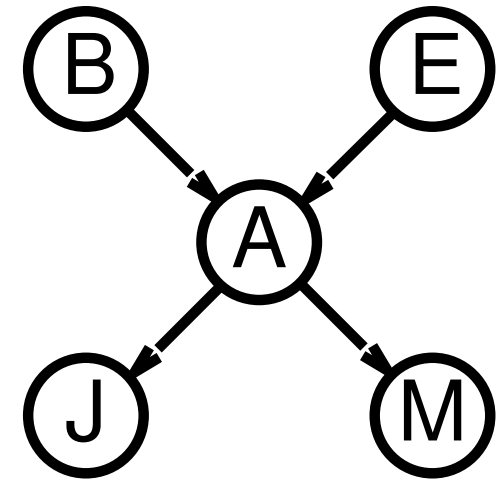


# Exemplo



# Representação compacta

- Cada linha numa CPT contém a probabilidade condicional de cada valor do nó para cada combinação de valores dos pais.
- Uma CPT para  $X_i$  Booleana com  $k$  pais Booleanos requer  $2^k$  linhas para as combinações de valores dos pais
- Cada linha requer um número  $p$  para  $X_i = \text{true}$  (o número para  $X_i = \text{false}$  é simplesmente  $1-p$ )
- Se cada variável não tem mais de  $k$  pais, a rede precisa no total de  $O(n \cdot 2^k)$  números
- I.e., cresce linearmente com  $n$ , vs.  $O(2^n)$  para a distribuição conjunta total
- Para a rede do assaltante,  $1+1+4+2+2=10$  números (vs.  $2^5-1 = 31$ )





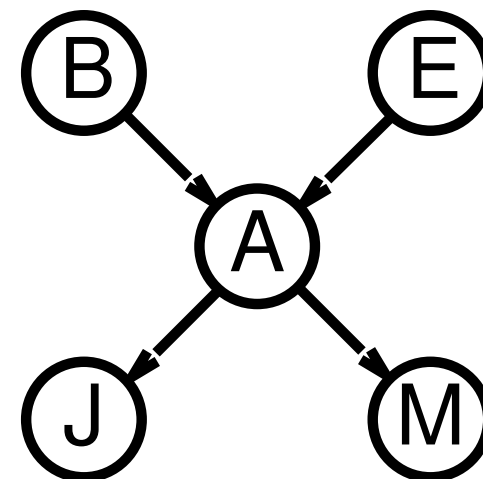
# Semântica

- Duas formas **equivalentes**:
  - **Semântica global** (ou numérica): interpretar as redes como uma representação da distribuição de probabilidade conjunta
    - Indica **como construir uma rede**
  - **Semântica local** (ou topológica): interpretar as redes como uma representação de uma coleção de declarações de independência condicional
    - Indica **como fazer inferências com uma rede**

# Semântica Global

- Semântica global define a distribuição conjunta total como o produto distribuições condicionais locais:

$$\mathbf{P}(X_1, \dots, X_n) = \prod_{i=1..n} \mathbf{P}(X_i \mid \text{Parents}(X_i))$$

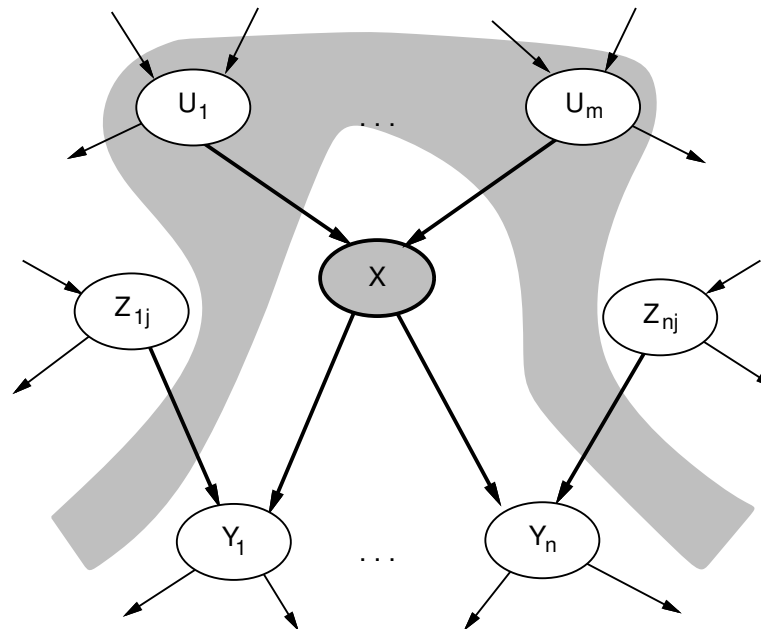


- Nesta situação dizemos que a distribuição  $\mathbf{P}$  é compatível com a rede  $G$ .
- E.g.

$$\begin{aligned} P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) &= \\ &= P(j|a) P(m|a) P(a|\neg b, \neg e) P(\neg b) P(\neg e) = \\ &= 0.90 \times 0.70 \times 0.001 \times 0.999 \times 0.998 = 0.00063 \end{aligned}$$

# Semântica Local

- Semântica local (topológica): cada nó é condicionalmente independente dos seus não-descendentes, dado os seus pais



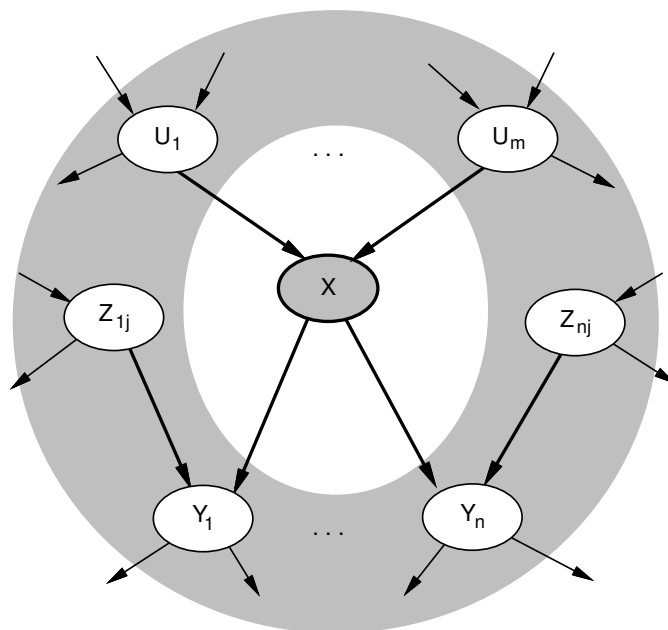
- Para todo o nó  $X$  assume-se que

$$P(X|Z_{1j}, \dots, Z_{nj}, U_1, \dots, U_m) = P(X|U_1, \dots, U_m)$$

- Conclui-se que: **Semântica Local  $\Leftrightarrow$  Semântica Global**

# Cobertura de Markov

- Cada nó é condicionalmente independente de todos os outros dada a sua cobertura de Markov (Markov blanket): pais + filhos + pais dos filhos



- Seja  $W_1, \dots, W_p$  um qualquer conjunto de nós da rede, disjunto da cobertura de Markov de X. Tem-se:

$$P(X|W_1, \dots, W_p, U_1, \dots, U_m, Y_1, \dots, Y_n, Z_{1j}, \dots, Z_{nj}) = P(X|U_1, \dots, U_m, Y_1, \dots, Y_n, Z_{1j}, \dots, Z_{nj})$$

# Construção de Redes Bayesianas

1. Escolher uma ordenação das variáveis  $X_1, \dots, X_n$
  2. Para  $i=1$  até  $n$ 
    1. Adicionar  $X_i$  à rede
    2. Escolher de entre  $X_1, \dots, X_{i-1}$  um conjunto minimal de pais para  $X_i$  tal que  $\mathbf{P}(X_i \mid \text{Parents}(X_i)) = \mathbf{P}(X_i \mid X_1, \dots, X_{i-1})$  e acrescentar ligações de cada pai para  $X_i$
    3. Escrever a tabela de probabilidade condicional (CPT) para  $\mathbf{P}(X_i \mid \text{Parents}(X_i))$
- Esta escolha de pais garante a semântica global:
  - $\mathbf{P}(X_1, \dots, X_n) = \prod_{i=1..n} \mathbf{P}(X_i \mid X_1, \dots, X_{i-1})$  (regra da cadeia)  
=  $\prod_{i=1..n} \mathbf{P}(X_i \mid \text{Parents}(X_i))$  (por construção)

# Exemplo

- Suponhamos a seguinte ordenação  $M, J, A, B, E$

MaryCalls

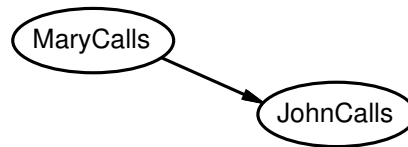
JohnCalls

Se a Mary tiver telefonado, isso provavelmente significa que o alarme disparou, o que tornaria mais provável que John também tivesse telefonado...

- $P(J|M) = P(J)$ ?

# Exemplo

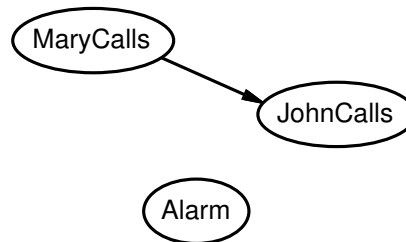
- Suponhamos a seguinte ordenação  $M, J, A, B, E$



- $P(J|M) = P(J)$ ? Não

# Exemplo

- Suponhamos a seguinte ordenação  $M, J, A, B, E$



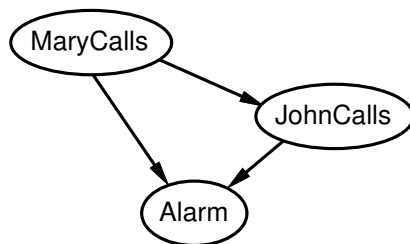
Claramente que se ambos tiverem telefonado, é mais provável que o alarme tenha disparado do que se apenas um deles tiver telefonado...

- $P(J|M) = P(J)$ ? Não
- $P(A|J,M) = P(A|J)$ ?  $P(A|J,M) = P(A)$ ?



# Exemplo

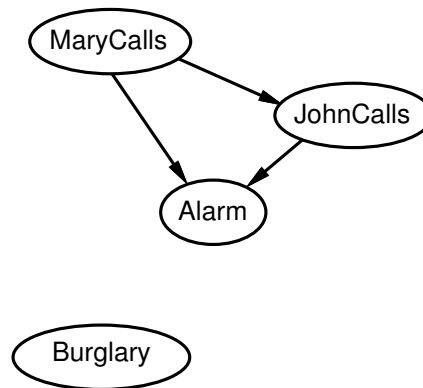
- Suponhamos a seguinte ordenação  $M, J, A, B, E$



- $P(J|M) = P(J)$ ? Não
- $P(A|J,M) = P(A|J)$ ?  $P(A|J,M) = P(A)$ ? Não

# Exemplo

- Suponhamos a seguinte ordenação  $M, J, A, B, E$

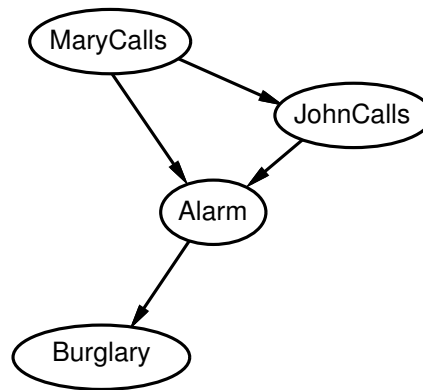


Se soubermos o estado do alarme, então os telefonemas não acrescentam nada em relação ao roubo...  
Apenas precisamos do alarme...

- $P(J|M) = P(J)$ ? Não
- $P(A|J,M) = P(A|J)$ ?  $P(A|J,M) = P(A)$ ? Não
- $P(B|A,J,M) = P(B|A)$ ?
- $P(B|A,J,M) = P(B)$ ?

# Exemplo

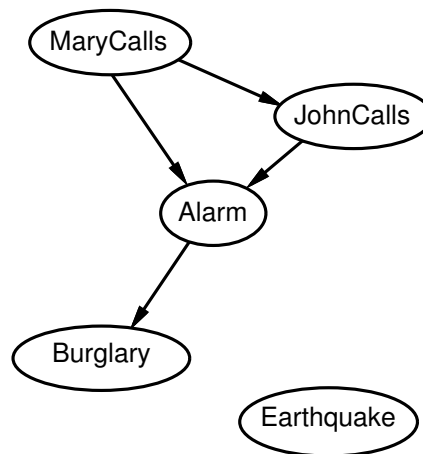
- Suponhamos a seguinte ordenação  $M, J, A, B, E$



- $P(J|M) = P(J)$ ? Não
- $P(A|J,M) = P(A|J)$ ?  $P(A|J,M) = P(A)$ ? Não
- $P(B|A,J,M) = P(B|A)$ ? Sim
- $P(B|A,J,M) = P(B)$ ? Não

# Exemplo

- Suponhamos a seguinte ordenação  $M, J, A, B, E$

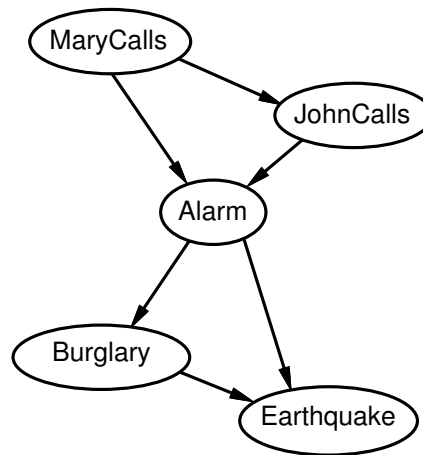


- $P(J|M) = P(J)$ ? Não
- $P(A|J,M) = P(A|J)$ ?  $P(A|J,M) = P(A)$ ? Não
- $P(B|A,J,M) = P(B|A)$ ? Sim
- $P(B|A,J,M) = P(B)$ ? Não
- $P(E|B,A,J,M) = P(E|A)$ ?
- $P(E|B,A,J,M) = P(E|A,B)$ ?

Se o alarme tiver disparado, é mais provável que tenha havido um tremor de terra. Mas se soubermos que houve um roubo, então isso explica o alarme, e a probabilidade de ter havido um tremor de terra é mais baixa...

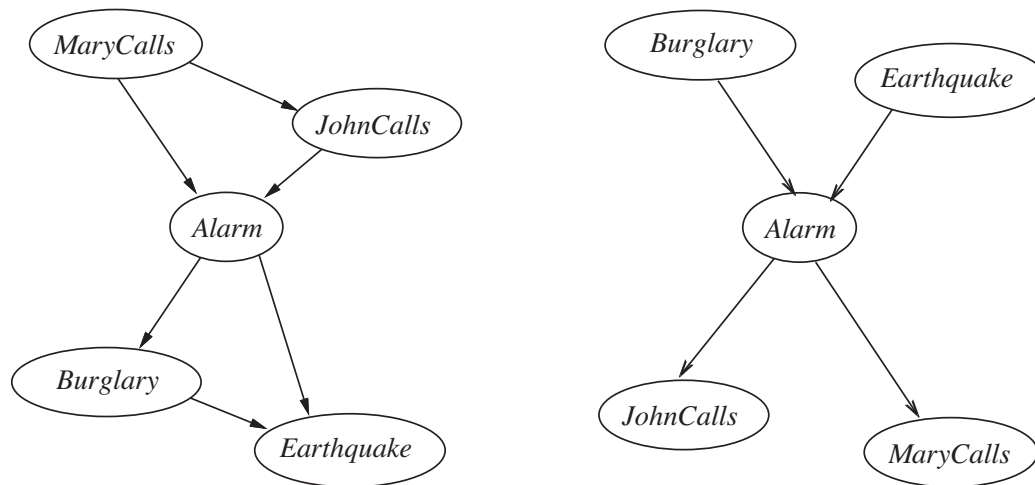
# Exemplo

- Suponhamos a seguinte ordenação  $M, J, A, B, E$



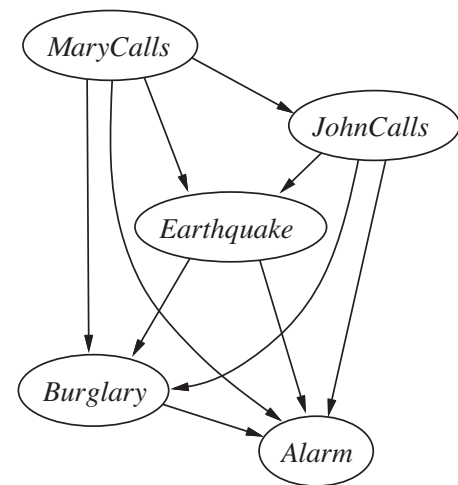
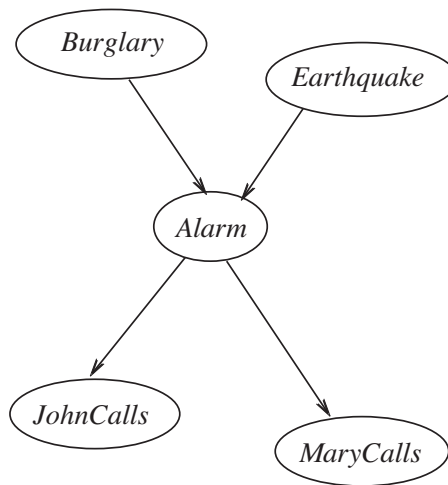
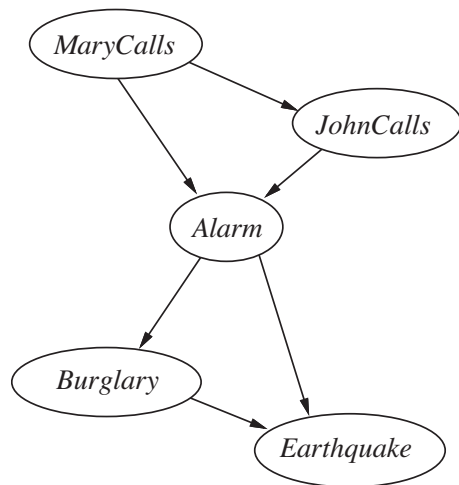
- $P(J|M) = P(J)$ ? Não
- $P(A|J,M) = P(A|J)$ ?  $P(A|J,M) = P(A)$ ? Não
- $P(B|A,J,M) = P(B|A)$ ? Sim
- $P(B|A,J,M) = P(B)$ ? Não
- $P(E|B,A,J,M) = P(E|A)$ ? Não
- $P(E|B,A,J,M) = P(E|A,B)$ ? Sim

# Exemplo



- Decidir sobre independência condicional é difícil nas direções não causais: modelos causais e independência condicional parecem ser inatos nos humanos!
- Determinação de probabilidades condicionais é difícil nas direções não causais.
  - Na rede Bayesiana da esquerda será necessário determinar a probabilidade condicional e.g. de *Earthquake* dado *Alarm* e *Burglary*, que não é fácil nem natural...
- Rede (esquerda) é menos compacta do que a inicial (direita). Necessários
  - Esquerda:  $1 + 2 + 4 + 2 + 4 = 13$  valores
  - Direita:  $1 + 1 + 4 + 2 + 2 = 10$  valores

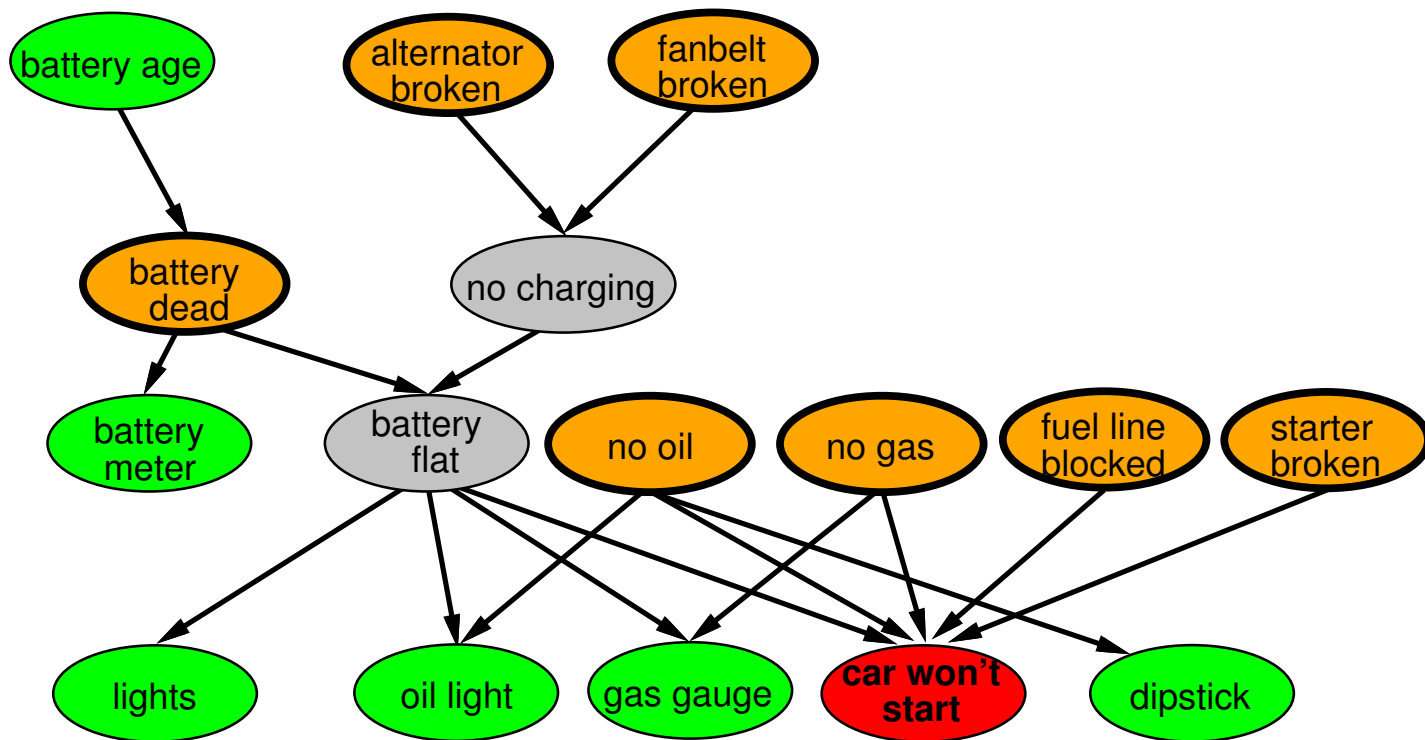
# Exemplo



- À direita, uma muito má ordenação das variáveis.
  - Requer a especificação de 31 valores (vs. 13 da esquerda e 10 da do meio)
- As três redes representam a mesma distribuição conjunta de probabilidade.
- No entanto, apenas a do meio expressa todas as independências condicionais, levando a uma representação mais compacta.

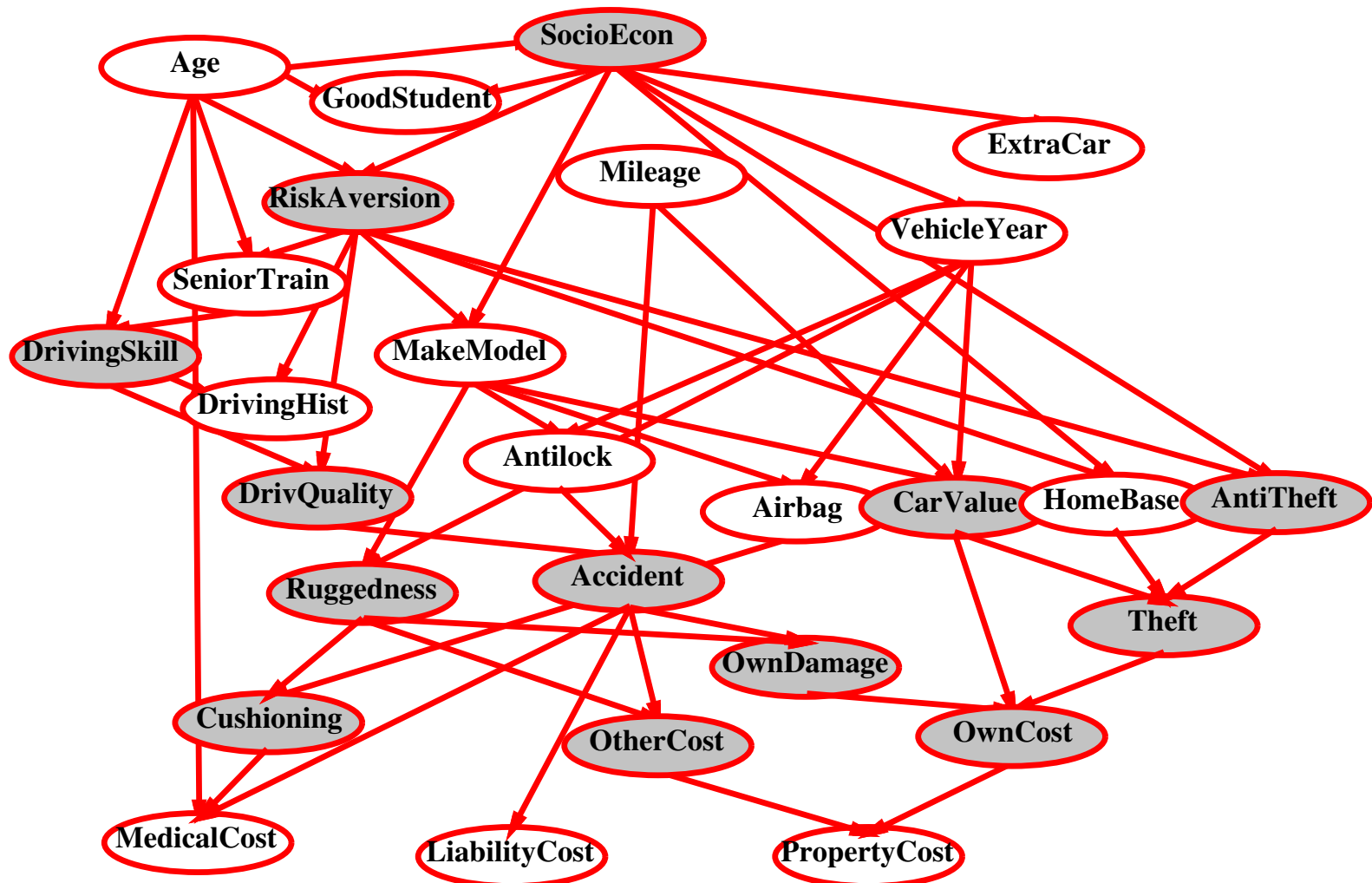
# Exemplo: diagnóstico de avarias

- Evidência inicial evidência: carro não pega
- Variáveis observáveis (verde), variáveis “avarias” (laranja)
- Variáveis ocultas (cinzento) garantem estrutura esparsa, reduzem parâmetros





# Exemplo: seguro do carro



# Distribuições condicionais compactas

- CPT cresce exponencialmente com o numero de pais
- CPT fica infinita com pai ou filho tomando valores contínuos
- Solução: distribuições canónicas que são definidas de forma compacta
- Nos determinísticos são o caso mais simples:
  - $X = f(\text{Parents}(X))$  para alguma função  $f$
- E.g., funções Booleanas
  - $\text{NorthAmerican} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$
- E.g., relações numéricas entre variáveis continuas
  - $\partial \text{Level} / \partial t = \text{inflow} + \text{precipitation} - \text{outflow} - \text{evaporation}$

# Distribuições condicionais compactas

- Distribuições **Noisy-OR** modelam múltiplas causas que não interagem
  - Pais  $U_1, \dots, U_k$  incluem todas as causas (pode-se adicionar nó)
  - Probabilidade de falha independentes  $q_i$  para cada causa isoladamente

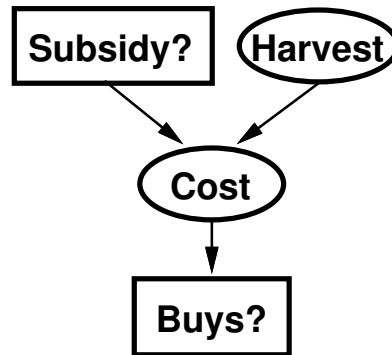
$$P(X|U_1, \dots, U_j, \neg U_{j+1}, \dots, \neg U_k) = 1 - \prod_{i=1..j} q_i$$

Constipação	Gripe	Malária	P(Febre)	P( $\neg$ Febre)
F	F	F	0.0	1.0
F	F	T	0.9	<b>0.1</b>
F	T	F	0.8	<b>0.2</b>
F	T	T	0.98	0.02=0.2×0.1
T	F	F	0.4	<b>0.6</b>
T	F	T	0.94	0.06=0.6×0.1
T	T	F	0.88	0.12=0.6×0.2
T	T	T	0.988	0.012=0.6×0.2×0.1

- Número de parâmetros linear no número de pais

# Redes Híbridas (discretas + contínuas)

- Discretas (*Subsidy?* e *Buys?*); contínuas (*Harvest* e *Cost*)



- Opção 1: discretização – erros podem ser grandes, CPTs grandes
- Opção 2: famílias canónicas com numero finito de parâmetros
  - Variável contínua, pais discretos+contínuos (e.g., *Cost*)
  - Variável discreta, pais contínuos (e.g., *Buys?*)

# Variáveis filho contínuas

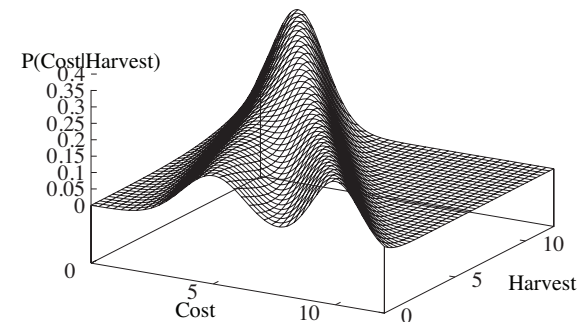
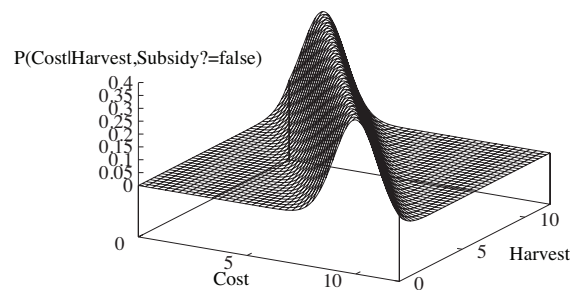
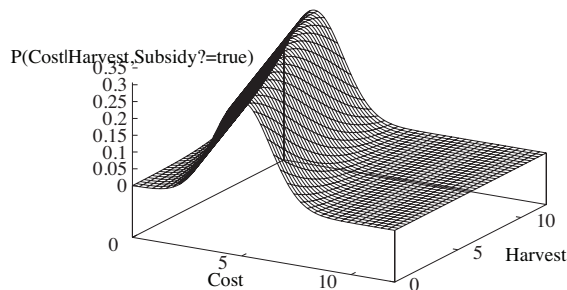
- Necessária uma função de **densidade condicional** para cada variável filho dados pais contínuos, para cada possível atribuição de pais discretos
- Mais habitual é o modelo **linear Gaussiano**, e.g.,:

$$\begin{aligned} &P(\textit{Cost} = c \mid \textit{Harvest} = h, \textit{Subsidy}? = \textit{true}) \\ &= N(a_t h + b_t, \sigma_t)(c) \\ &= \frac{1}{\sigma_t \sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{c - (a_t h + b_t)}{\sigma_t}\right)^2\right) \end{aligned}$$

- Média **Cost** varia linearmente com **Harvest**, variância fixa
- Variação linear não é razoável para todo o domínio mas funciona bem se os valores **prováveis** de **Harvest** estiverem limitados

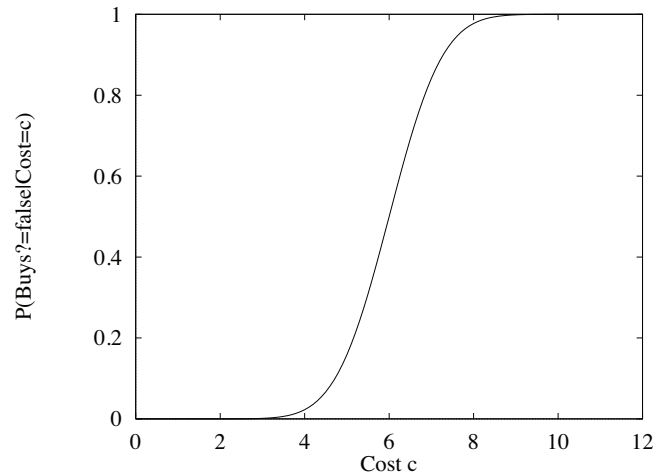
# Variáveis filho contínuas

- Rede só com variáveis contínuas com distribuição Linear Gaussiana (LG)
  - distribuição conjunta tem distribuição Gaussiana multivariada
- Rede discreta+contínua LG é uma rede **Gaussiana condicional**, i.e., uma Gaussiana multivariada para todas as variáveis contínuas, para cada combinação de valores de variáveis discretas



# Variáveis discretas com pais contínuos

- Probabilidade de *Buys?* dado *Cost* deve ser um limiar “suave”:



- **Distribuição Probit** utiliza integral de uma Gaussiana:

$$\Phi(x) = \int_{-\infty}^x N(0,1)(x)$$

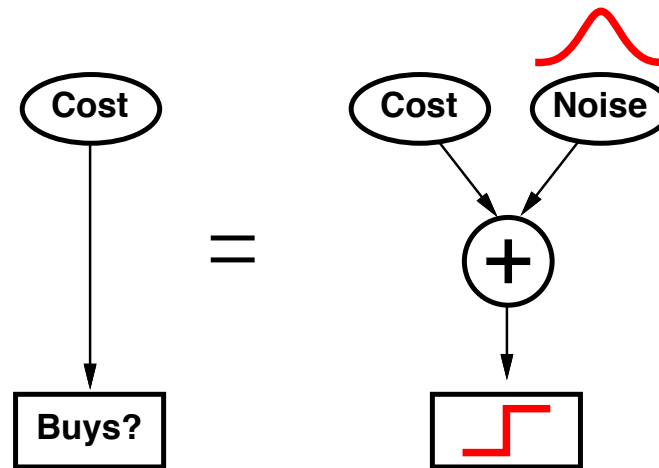
$$P(Buys? = true | Cost = c) = \Phi((-c + \mu) / \sigma)$$

$$P(Buys? = false | Cost = c) = 1 - \Phi((-c + \mu) / \sigma) = \Phi((c - \mu) / \sigma)$$

- Em que  $\mu$  é o local onde ocorre o limiar e  $\sigma$  um parâmetro que controla a largura do limiar.

# Porquê a probit

1. Tem a forma correcta
2. Pode ser entendida como um limiar cuja localização esta sujeita a ruído



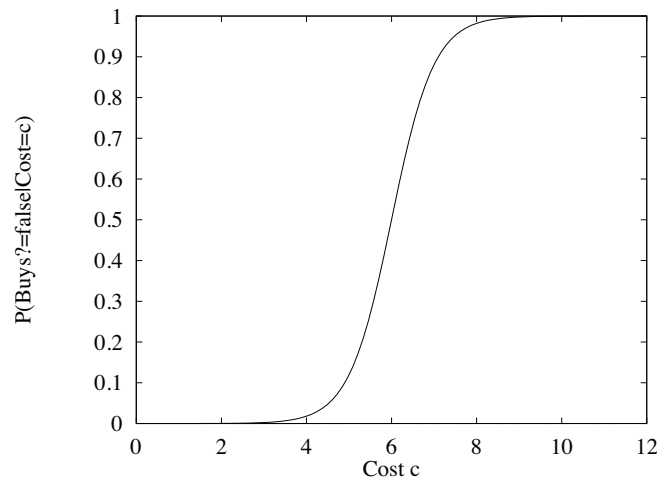


# Variável discreta (cont.)

- Distribuição Sigmoidal (ou logit) também utilizada em redes neurais:

$$P(Buys? = true \mid Cost = c) = \frac{1}{1 + \exp\left(-2 \frac{-c + \mu}{\sigma}\right)}$$

- Sigmoidal tem forma semelhante a da probit mas com caudas maiores



# Sumário

- Redes de Bayes são uma representação natural para independência condicional (induzidas causalmente)
- Topologia + CPTs = representação compacta de distribuição conjunta
- Geralmente fácil de construir por (não)peritos
- Distribuições canônicas (e.g., noisy-OR) = representação compacta de CPTs
- Variáveis contínuas  $\Rightarrow$  distribuições parametrizadas (e.g., Gaussiana linear)