# KI2 - 4

# Bayesian Belief Networks

*AIMA, Chapter 14*

Kunstmatige Intelligentie / RuG

Initial evidence: car won't start
Testable variables (green), "broken, so fix it" variables (orange)
Hidden variables (gray) ensure sparse structure, reduce parameters
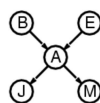


---

# Bayesian Networks

- A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions

- The constituents of a Bayesian network:
  - a set of nodes, one per variable
  - a directed, acyclic graph (link ≈ "directly influences")
  - a conditional distribution for each node given its parents
    $P(X_i \mid$ Parents $(X_i))$

- In the simplest case, conditional distribution represented as a conditional probability table (CPT) giving the distribution over $X_i$ for each combination of parent values

---

# Example

- Topology of network encodes conditional independence assertions:



- *Weather* is independent of the other variables
- *Toothache* and *Catch* are conditionally independent given *Cavity* (e.g.: p(T|Cav) = p(T|Cav, Catch)
  =p(T|Cav, not Catch)

---

# Example

- I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?

- Variables: *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

- Network topology reflects "causal" knowledge:
  - A burglar can set the alarm off
  - An earthquake can set the alarm off
  - The alarm can cause Mary to call
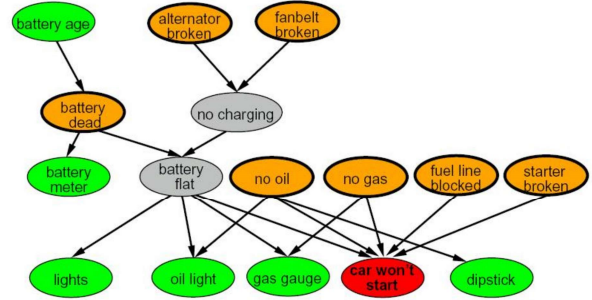  - The alarm can cause John to call

---

# Example contd.



---

# Compactness

- A CPT for Boolean $X_i$ with $k$ Boolean parents has $2^k$ rows for the combinations of parent values
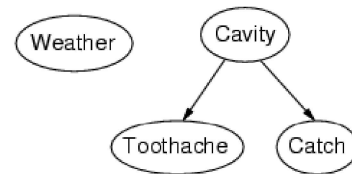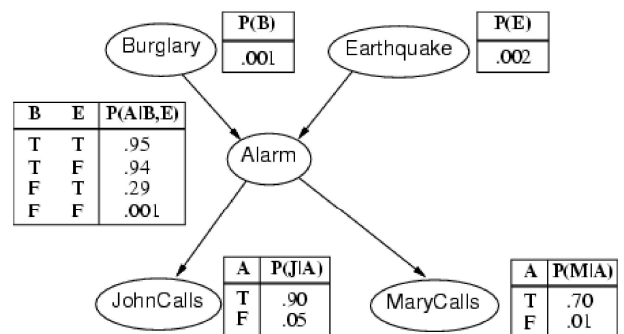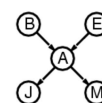


- Each row requires one number $p$ for $X_i$ = *true* (the number for $X_i$ = *false* is just *1-p*)
- If each variable has no more than $k$ parents, the complete network requires $O(n \cdot 2^k)$ numbers
  i.e. grows linearly with $n$, vs. $O(2^n)$ for the full joint distribution
- For burglary net, 1 + 1 + 4 + 2 + 2 = 10 numbers (vs. $2^5$-1 = 31)

---

# Computing the Full Joint Distribution



- The full joint distribution is equal to the product of the local conditional distributions:

$$P(X_1, \ldots, X_n) = \Pi_i \, P(X_i \mid Parents(X_i))$$

e.g. $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$
$= P(j \mid a) \, P(m \mid a) \, P(a \mid \neg b, \neg e) \, P(\neg b) \, P(\neg e)$

# Constructing Bayesian Networks

1. Choose an ordering of variables $X_1, \ldots, X_n$
2. For $i = 1$ to $n$
   - add $X_i$ to the network
   - select parents from $X_1, \ldots, X_{i-1}$ such that
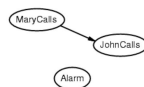     $P(X_i \mid Parents(X_i)) = P(X_i \mid X_1, \ldots X_{i-1})$

- This choice of parents guarantees:

$P(X_1, \ldots, X_n) = \prod_i P(X_i \mid X_1, \ldots, X_{i-1})$ (chain rule)

$\qquad\qquad = \prod_i P(X_i \mid Parents(X_i))$ (by construction)

# Example

- Suppose we choose the ordering *M, J, A, B, E*

MaryCalls

JohnCalls

$P(J \mid M) = P(J)?$

# Example

- Suppose we choose the ordering *M, J, A, B, E*

MaryCalls → JohnCalls

Alarm

$P(J \mid M) = P(J)?$ No, John can decide for himself
$P(A \mid J, M) = P(A \mid J)?\ P(A \mid J, M) = P(A)?$

# Example

- Suppose we choose the ordering *M, J, A, B, E*

MaryCalls → JohnCalls

Alarm

Burglary

$P(J \mid M) = P(J)?$ No
$P(A \mid J, M) = P(A \mid J)?\ P(A \mid J, M) = P(A)?$ No
$P(B \mid A, J, M) = P(B \mid A)?$
$P(B \mid A, J, M) = P(B)?$

# Example

- Suppose we choose the ordering *M, J, A, B, E*

MaryCalls → JohnCalls

Alarm

Burglary

Earthquake

$P(J \mid M) = P(J)?$ No
$P(A \mid J, M) = P(A \mid J)?\ P(A \mid J, M) = P(A)?$ No
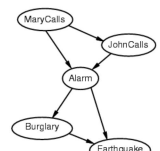$P(B \mid A, J, M) = P(B \mid A)?$ Yes
$P(B \mid A, J, M) = P(B)?$ No, the burglar decides himself
$P(E \mid B, A, J, M) = P(E \mid A)?$
$P(E \mid B, A, J, M) = P(E \mid A, B)?$

# Example

- Suppose we choose the ordering *M, J, A, B, E*

MaryCalls → JohnCalls

Alarm

Burglary

Earthquake

$P(J \mid M) = P(J)?$ No
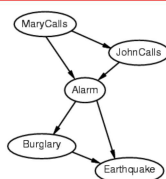$P(A \mid J, M) = P(A \mid J)?\ P(A \mid J, M) = P(A)?$ No
$P(B \mid A, J, M) = P(B \mid A)?$ Yes
$P(B \mid A, J, M) = P(B)?$ No, the burglar decides himself
$P(E \mid B, A, J, M) = P(E \mid A)?$ No, earthquakes don't wait for an alarm to go off
$P(E \mid B, A, J, M) = P(E \mid A, B)?$ Yes  ➔ *unclear picture*

# Example contd.
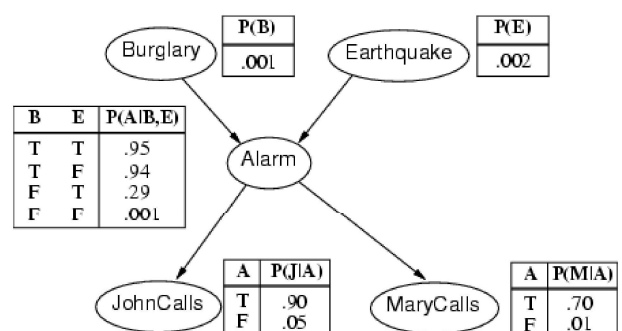
MaryCalls → JohnCalls

Alarm

Burglary

Earthquake

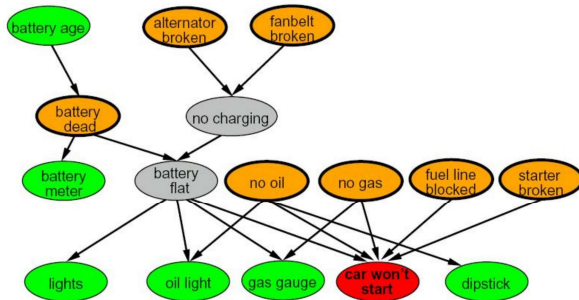- Deciding conditional independence is difficult in **noncausal** directions.
- Causal models and conditional independence seem hardwired in human reasoning!
- Also: the bad network was less compact:
  1 + 2 + 4 + 2 + 4 = 13 numbers needed.

# Causal ordering is natural, easy

| P(B) |
|------|
| .001 |

| P(E) |
|------|
| .002 |

| B | E | P(A\|B,E) |
|---|---|-----------|
| T | T | .95 |
| T | F | .94 |
| F | T | .29 |
| F | F | .001 |

Burglary → Alarm ← Earthquake

Alarm → JohnCalls

Alarm → MaryCalls

| A | P(J\|A) |
|---|---------|
| T | .90 |
| F | .05 |

| A | P(M\|A) |
|---|---------|
| T | .70 |
| F | .01 |

Initial evidence: car won't start
Testable variables (green), "broken, so fix it" variables (orange)
Hidden variables (gray) ensure sparse structure, reduce parameters

# Summary

- Bayesian networks provide a natural representation for (causally induced) conditional independence.
- Topology + CPTs = compact representation of joint distribution.
- Generally easy for domain experts to construct.
- Belief networks have found increasing use in complex diagnosis problems (medical, cars, PC operating systems).

## All is well in Belief Network Land?

- Problems:
  – Network construction: often gofai human construction labor (knowledge based, 'Cyc' etc.)
  – Estimation of probabilities
  – Product of probabilities:
    • Not p but p' = p + ε   therefore

$$\Pi_i \; (p_i + \varepsilon)$$ propagation of errors! especially in a long chain

## Summary for Bayesian Methods

- Bayesian methods:
  - Learning = estimation of probability distributions of samples from different classes
  - Classification = use these estimates to determine which class is more likely for a new instance
- Naive Bayes:
  - Assumes that attributes are independent.
- Bayesian Belief Networks:
  - Assumes that subsets of attributes are independent.
- Bayesian methods allow combining prior knowledge about the world with evidence from the data stream.