## CS 343: Artificial Intelligence
## Bayesian Networks

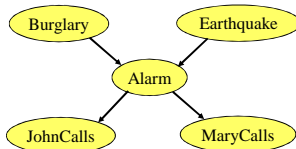Raymond J. Mooney

University of Texas at Austin

1

---

## Graphical Models

- If no assumption of independence is made, then an exponential number of parameters must be estimated for sound probabilistic inference.
- No realistic amount of training data is sufficient to estimate so many parameters.
- If a blanket assumption of conditional independence is made, efficient training and inference is possible, but such a strong assumption is rarely warranted.
- **Graphical models** use directed or undirected graphs over a set of random variables to explicitly specify variable dependencies and allow for less restrictive independence assumptions while limiting the number of parameters that must be estimated.
  - **Bayesian Networks**: Directed acyclic graphs that indicate causal structure.
  - **Markov Networks**: Undirected graphs that capture general dependencies.
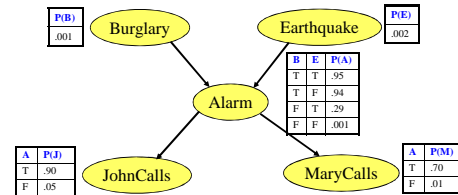
2

---

## Bayesian Networks

- Directed Acyclic Graph (DAG)
  - Nodes are random variables
  - Edges indicate causal influences



3

---

## Conditional Probability Tables

- Each node has a **conditional probability table** (**CPT**) that gives the probability of each of its values given every possible combination of values for its parents (conditioning case).
  - Roots (sources) of the DAG that have no parents are given prior probabilities.



4

---

## CPT Comments

- Probability of false not given since rows must add to 1.
- Example requires 10 parameters rather than $2^5 - 1 = 31$ for specifying the full joint distribution.
- Number of parameters in the CPT for a node is exponential in the number of parents (fan-in).

5

---

## Joint Distributions for Bayes Nets

- A Bayesian Network implicitly defines a joint distribution.

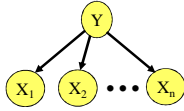$$P(x_1, x_2, ... x_n) = \prod_{i=1}^{n} P(x_i \mid \text{Parents}(X_i))$$

- Example

$$P(J \wedge M \wedge A \wedge \neg B \wedge \neg E)$$
$$= P(J \mid A)P(M \mid A)P(A \mid \neg B \wedge \neg E)P(\neg B)P(\neg E)$$
$$= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 = 0.00062$$

- Therefore an inefficient approach to inference is:
  - 1) Compute the joint distribution using this equation.
  - 2) Compute any desired conditional probability using the joint distribution.

6

1

## Naïve Bayes as a Bayes Net

- Naïve Bayes is a simple Bayes Net



- Priors P($Y$) and conditionals P($X_i|Y$) for Naïve Bayes provide CPTs for the network.

7

## Independencies in Bayes Nets

- If removing a subset of nodes $S$ from the network renders nodes $X_i$ and $X_j$ disconnected, then $X_i$ and $X_j$ are independent given $S$, i.e. P($X_i | X_j$, $S$) = P($X_i | S$)
- However, this is too strict a criteria for conditional independence since two nodes will still be considered independent if their simply exists some variable that depends on both.
  – For example, Burglary and Earthquake should be considered independent since they both cause Alarm.
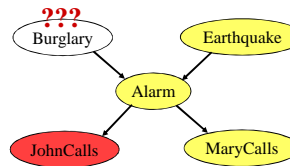
8

## Independencies in Bayes Nets (cont.)

- Unless we know something about a common effect of two "independent causes" or a descendent of a common effect, then they can be considered independent.
  – For example, if we know nothing else, Earthquake and Burglary are independent.
- However, if we have information about a common effect (or descendent thereof) then the two "independent" causes become probabilistically linked since evidence for one cause can "explain away" the other.
  – For example, if we know the alarm went off that someone called about the alarm, then it makes earthquake and burglary dependent since evidence for earthquake decreases belief in burglary. and vice versa.

9

## Bayes Net Inference

- Given known values for some **evidence variables**, determine the posterior probability of some **query variables**.
- Example: Given that John calls, what is the probability that there is a Burglary?
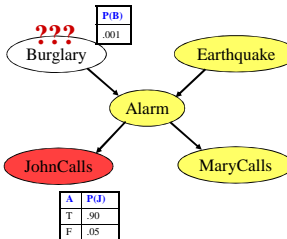


John calls 90% of the time there is an Alarm and the Alarm detects 94% of Burglaries so people generally think it should be fairly high.

However, this ignores the prior probability of John calling.

10

## Bayes Net Inference

- Example: Given that John calls, what is the probability that there is a Burglary?
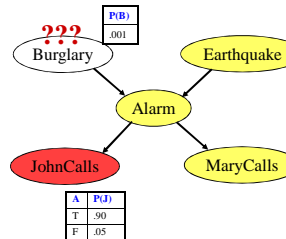


John also calls 5% of the time when there is no Alarm. So over 1,000 days we expect 1 Burglary and John will probably call. However, he will also call with a false report 50 times on average. So the call is about 50 times more likely a false report: P(Burglary | JohnCalls) ≈ 0.02
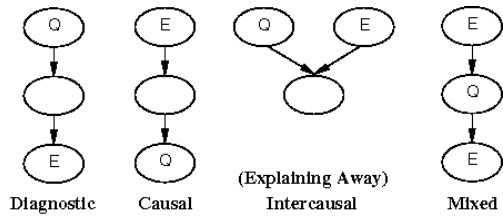
11

## Bayes Net Inference

- Example: Given that John calls, what is the probability that there is a Burglary?



Actual probability of Burglary is 0.016 since the alarm is not perfect (an Earthquake could have set it off or it could have gone off on its own). On the other side, even if there was not an alarm and John called incorrectly, there could have been an undetected Burglary anyway, but this is unlikely.

12

2

## Types of Inference



Diagnostic    Causal    (Explaining Away) Intercausal    Mixed

## Sample Inferences

- **Diagnostic (evidential, abductive)**: From effect to cause.
  - P(Burglary | JohnCalls) = 0.016
  - P(Burglary | JohnCalls ∧ MaryCalls) = 0.29
  - P(Alarm | JohnCalls ∧ MaryCalls) = 0.76
  - P(Earthquake | JohnCalls ∧ MaryCalls) = 0.18
- **Causal (predictive)**: From cause to effect
  - P(JohnCalls | Burglary) = 0.86
  - P(MaryCalls | Burglary) = 0.67
- **Intercausal (explaining away)**: Between causes of a common effect.
  - P(Burglary | Alarm) = 0.376
  - P(Burglary | Alarm ∧ Earthquake) = 0.003
- **Mixed**: Two or more of the above combined
  - (diagnostic and causal) P(Alarm | JohnCalls ∧ ¬Earthquake) = 0.03
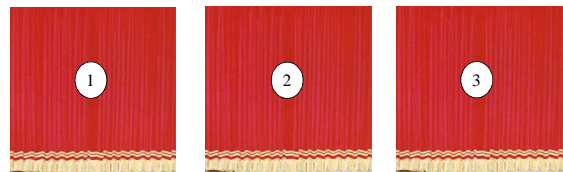  - (diagnostic and intercausal) P(Burglary | JohnCalls ∧ ¬Earthquake) = 0.017

## Probabilistic Inference in Humans

- People are notoriously bad at doing correct probabilistic reasoning in certain cases.
- One problem is they tend to ignore the influence of the prior probability of a situation.

## Monty Hall Problem



**One Line Demo:**
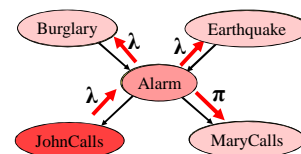http://math.ucsd.edu/~crypto/Monty/monty.html

## Complexity of Bayes Net Inference

- In general, the problem of Bayes Net inference is NP-hard (exponential in the size of the graph).
- For **singly-connected networks** or **polytrees** in which there are no undirected loops, there are linear-time algorithms based on **belief propagation**.
  - Each node sends local evidence messages to their children and parents.
  - Each node updates belief in each of its possible values based on incoming messages from it neighbors and propagates evidence on to its neighbors.
- There are approximations to inference for general networks based on **loopy belief propagation** that iteratively refines probabilities that converge to accurate values in the limit.

## Belief Propagation Example

- λ messages are sent from children to parents representing abductive evidence for a node.
- π messages are sent from parents to children representing causal evidence for a node.

## Belief Propagation Details

- Each node $B$ acts as a simple processor which maintains a vector $\lambda(B)$ for the total evidential support for each value of its corresponding variable and an analogous vector $\pi(B)$ for the total causal support.
- The belief vector $BEL(B)$ for a node, which maintains the probability for each value, is calculated as the normalized product:

$$BEL(B) = \alpha \, \lambda(B) \, \pi(B)$$

- Computation at each node involve $\lambda$ and $\pi$ **message vectors** sent between nodes and consists of simple matrix calculations using the CPT to update belief (the $\lambda$ and $\pi$ node vectors) for each node based on new evidence.

19

## Belief Propagation Details (cont.)

- Assumes the CPT for each node is a matrix ($M$) with a column for each value of the node's variable and a row for each conditioning case (all rows must sum to 1).
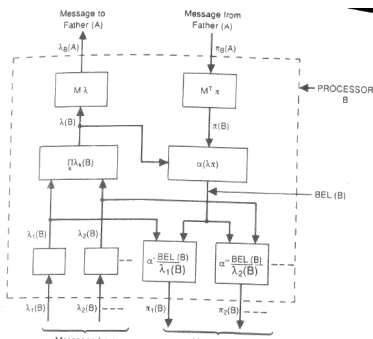
|  | Value of Alarm | |
|---|---|---|
|  | T | F |
| Values TT | 0.95 | 0.05 |
| of Burglary TF | 0.94 | 0.06 |
| and Earthquake FT | 0.29 | 0.71 |
| FF | 0.001 | 0.999 |

Matrix $M$ for the Alarm node

- Propagation algorithm is simplest for trees in which each node has only one parent (i.e. one cause).
- To initialize, $\lambda(B)$ for all leaf nodes is set to all 1's and $\pi(B)$ of all root nodes is set to the priors given in the CPT. Belief based on the root priors is then propagated down the tree to all leaves to establish priors for all nodes.
- Evidence is then added incrementally and the effects propagated to other nodes.
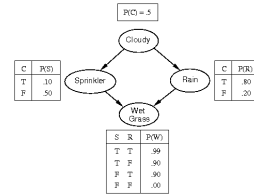
20

## Processor for Tree Networks



21

## Multiply Connected Networks

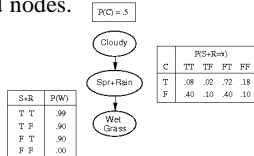- Networks with undirected loops, more than one directed path between some pair of nodes.



- In general, inference in such networks is NP-hard.
- Some methods construct a polytree(s) from given network and perform inference on transformed graph.

22

## Node Clustering

- Eliminate all loops by merging nodes to create **meganodes** that have the cross-product of values of the merged nodes.



- Number of values for merged node is exponential in the number of nodes merged.
- Still reasonably tractable for many network topologies requiring relatively little merging to eliminate loops.

23

## Bayes Nets Applications

- Medical diagnosis
  - Pathfinder system outperforms leading experts in diagnosis of lymph-node disease.
- Microsoft applications
  - Problem diagnosis: printer problems
  - Recognizing user intents for HCI
- Text categorization and spam filtering
- Student modeling for intelligent tutoring systems.

24

4

## Statistical Revolution

- Across AI there has been a movement from logic-based approaches to approaches based on probability and statistics.
  – Statistical natural language processing
  – Statistical computer vision
  – Statistical robot navigation
  – Statistical learning
- Most approaches are feature-based and "propositional" and do not handle complex relational descriptions with multiple entities like those typically requiring predicate logic.
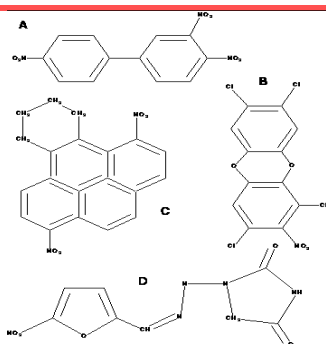
25

## Structured (Multi-Relational) Data

- In many domains, data consists of an unbounded number of entities with an arbitrary number of properties and relations between them.
  – Social networks
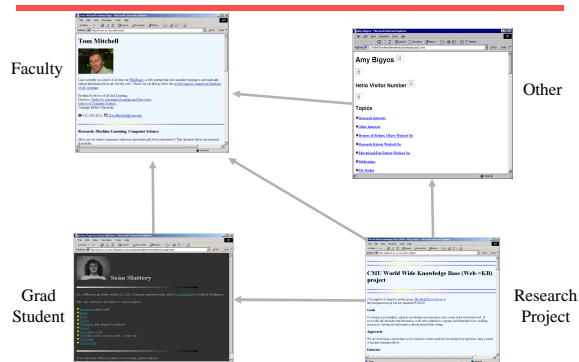  – Biochemical compounds
  – Web sites

## Biochemical Data



**Predicting mutagenicity**
[Srinivasan et. al, 1995]

27

## Web-KB Dataset [Slattery & Craven, 1998]



Faculty

Other

Grad Student

Research Project

## Collective Classification

- Traditional learning methods assume that objects to be classified are independent (the first "i" in the i.i.d. assumption)
- In structured data, the class of an entity can be influenced by the classes of related entities.
- Need to assign classes to all objects simultaneously to produce the most probable globally-consistent interpretation.

## Logical AI Paradigm

- Represents knowledge and data in a binary symbolic logic such as FOPC.
- + Rich representation that handles arbitrary sets of objects, with properties, relations, quantifiers, etc.
- − Unable to handle uncertain knowledge and probabilistic reasoning.

## Probabilistic AI Paradigm

- Represents knowledge and data as a fixed set of random variables with a joint probability distribution.
- + Handles uncertain knowledge and probabilistic reasoning.
- – Unable to handle arbitrary sets of objects, with properties, relations, quantifiers, etc.

## Statistical Relational Models

- Integrate methods from predicate logic (or relational databases) and probabilistic graphical models to handle structured, multi-relational data.
  - Probabilistic Relational Models (PRMs)
  - Stochastic Logic Programs (SLPs)
  - Bayesian Logic Programs (BLPs)
  - Relational Markov Networks (RMNs)
  - Markov Logic Networks (MLNs)
  - Other TLAs

32

## Conclusions

- Bayesian learning methods are firmly based on probability theory and exploit advanced methods developed in statistics.
- Naïve Bayes is a simple generative model that works fairly well in practice.
- A Bayesian network allows specifying a limited set of dependencies using a directed graph.
- Inference algorithms allow determining the probability of values for query variables given values for evidence variables.

33