# Bayesian Classifiers

CS510

Lecture #18

4/22/13

# Where are we?

- Learning the basics of classifiers
  - Goal: become an intelligent user
  - SVMs : Linear classifiers
    - Learn more in CS548
  - Feedforward Networks : Non-linear classifiers
    - Learn more in CS545
  - Today: Bayesian classifiers
    - Learn more in CS440

**Colorado State University**

# Review: Probability Basics

- Let X be variable whose value takes on a discrete set of labels, X = {x$_1$, …., x$_n$}
- The P(X=x$_i$) is a probability function if:

$$\forall x_i : 0 \le p(X = x_i) \le 1$$

$$\sum_i p(X = x_i) = 1$$

- We abbreviate P(X=x$_i$) as P(x$_i$)

Colorado State University

# Review: Probability Basics (II)

- Probabilities reflect the likelihood that a statement $X = x_i$ is true
  - If the statement is true, $P(x_i) = 1$
  - If it is false, $P(x_i) = 0$
  - Otherwise, higher values indicate more likely
- *Frequentist* probabilities represent samplings of random draws
- Subjective probabilities may not

Colorado State University

# Bayes Rule

- Given $X = \left\{ x_1, \ldots, x_n \right\}, Y = \left\{ y_1, \ldots, y_m \right\}$

$$P\left( x_i \wedge y_j \right) = P\left( x_i \mid y_j \right) P\left( y_j \right) = P\left( y_j \mid x_i \right) P\left( x_i \right)$$

- Or, put another way,

$$P\left( x_i \mid y_j \right) = \frac{P\left( y_j \mid x_i \right) P\left( x_i \right)}{P\left( y_j \right)}$$

# Bayesian Classification Example

- Let's say you extract a circle from an image, and want to know if it's a tire?

$$\text{Object(O)} = \{Wheel, handlebar, road, dirt, \ldots\}$$

$$\text{Feature(F)} = \{Circle, \neg Circle\}$$

$$P(Wheel \mid Circle\} = \frac{P(Circle \mid Wheel)P(Wheel)}{P(Circle)}$$

- Note that P(Circle|Wheel) is easier to estimate than P(Wheel|Circle)

# Naïve Bayes Classifiers

- Assume that features are independent, so

$$P(x \mid f_1, \ldots, f_m) = P(x \mid f_1) P(x \mid f_2) \ldots P(x \mid f_m)$$

- And of course

$$P(x \mid f_i) = \frac{P(f_i \mid x) P(x)}{P(f_i)}$$

- So

$$P(x \mid f_1, \ldots, f_m) = P(f_1 \mid x) \ldots P(f_m \mid x) \frac{P(x)^m}{\Pi P(f_i)}$$

# Naïve Bayes Classifiers (II)

- The fractional term in the last equation is constant for all $x_i$.

- So the most likely x is the one that maximizes

$$P(f_1 \mid x) \ldots P(f_m \mid x)$$

- You can recover the true probabilities by normalizing for all $x_i$

- Works well with PCA (where features are approximately independent)

# Conditional Independence

- Unfortunately, most random variables of interest are not independent

- A more useful notion is *conditional independence*

- Two variables X and Y are conditionally independent given Z if

  – $P(X = x | Y = y, Z=z) = P(X = x | Z=z)$ for all values x,y,z

  – That is, learning the values of Y does not change prediction of X once we know the value of Z

  – Notation: $I( X ; Y | Z )$

# Conditionally Chained Inference

- Independence is a strong assumption
- Often, x depends on multiple features that are not independent of each other
- If the features can be chained so that each depends only on previous features other…

$$P\left(x \mid f_1, f_2, f_3\right) = P\left(x \mid f_1\right) P\left(f_1 \mid f_2, f_3\right) P\left(f_2, f_3\right)$$

$$= P\left(x \mid f_1\right) P\left(f_1 \mid f_2\right) P\left(f_2 \mid f_3\right) P\left(f_3\right)$$

# Purpose of Bayesian Networks

- Facilitate the description of a collection of beliefs by
  - making causality relations explicit
  - exploiting conditional independence
- Provide efficient methods for:
  - Representing a joint probability distribution
  - Updating belief strengths when new evidence is observed

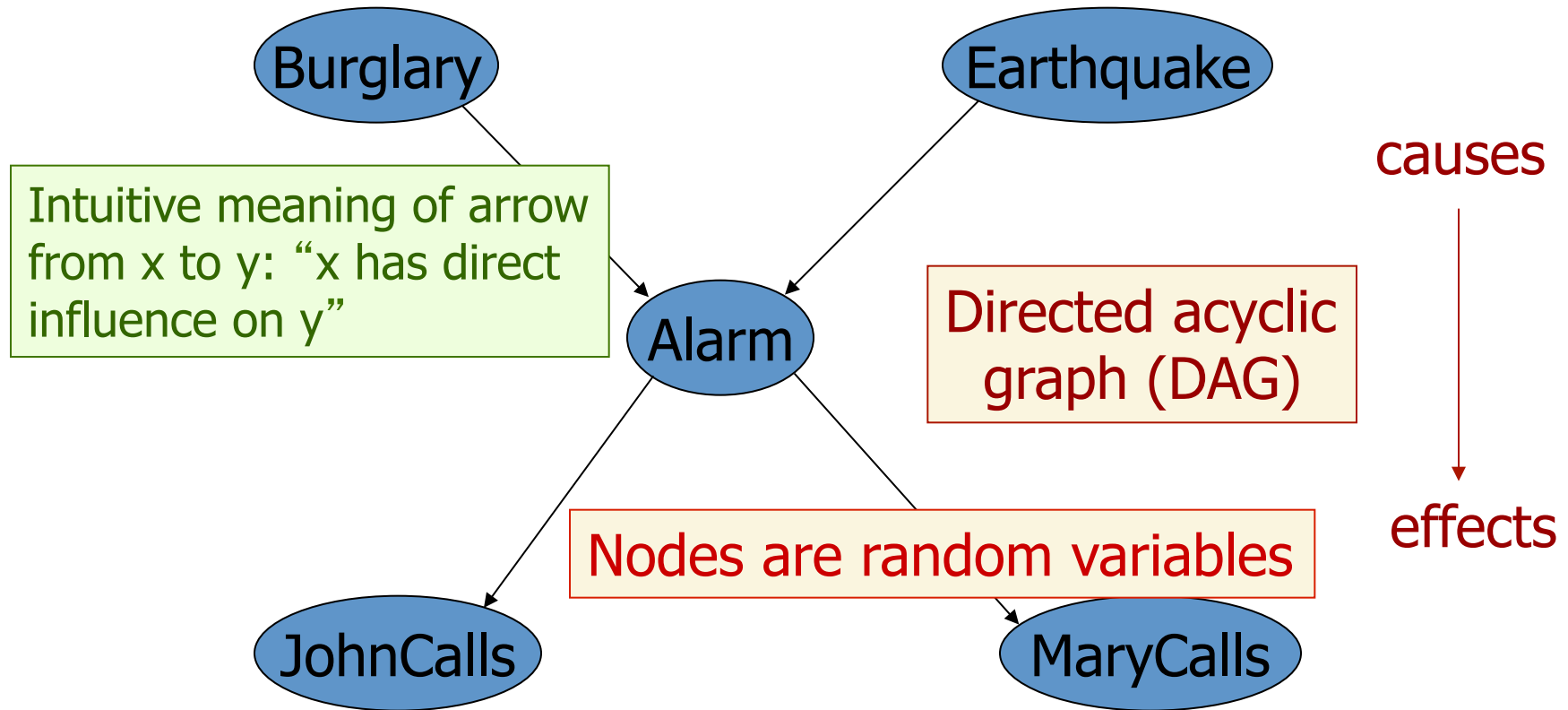**Colorado State University**

# Example

*I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometime it's set off by a minor earthquake. Is there a burglary?*

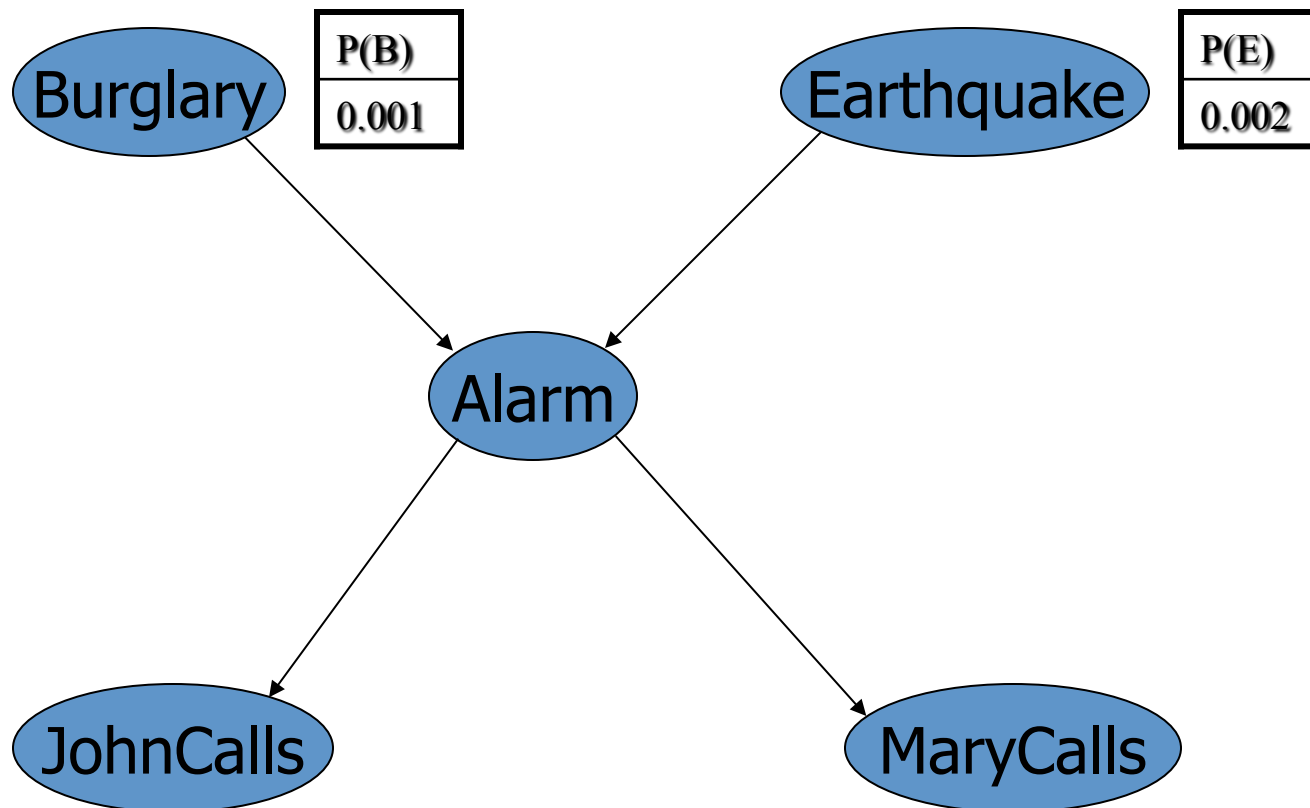Variables: Burglary, Earthquake, Alarm, JohnCalls, MaryCalls

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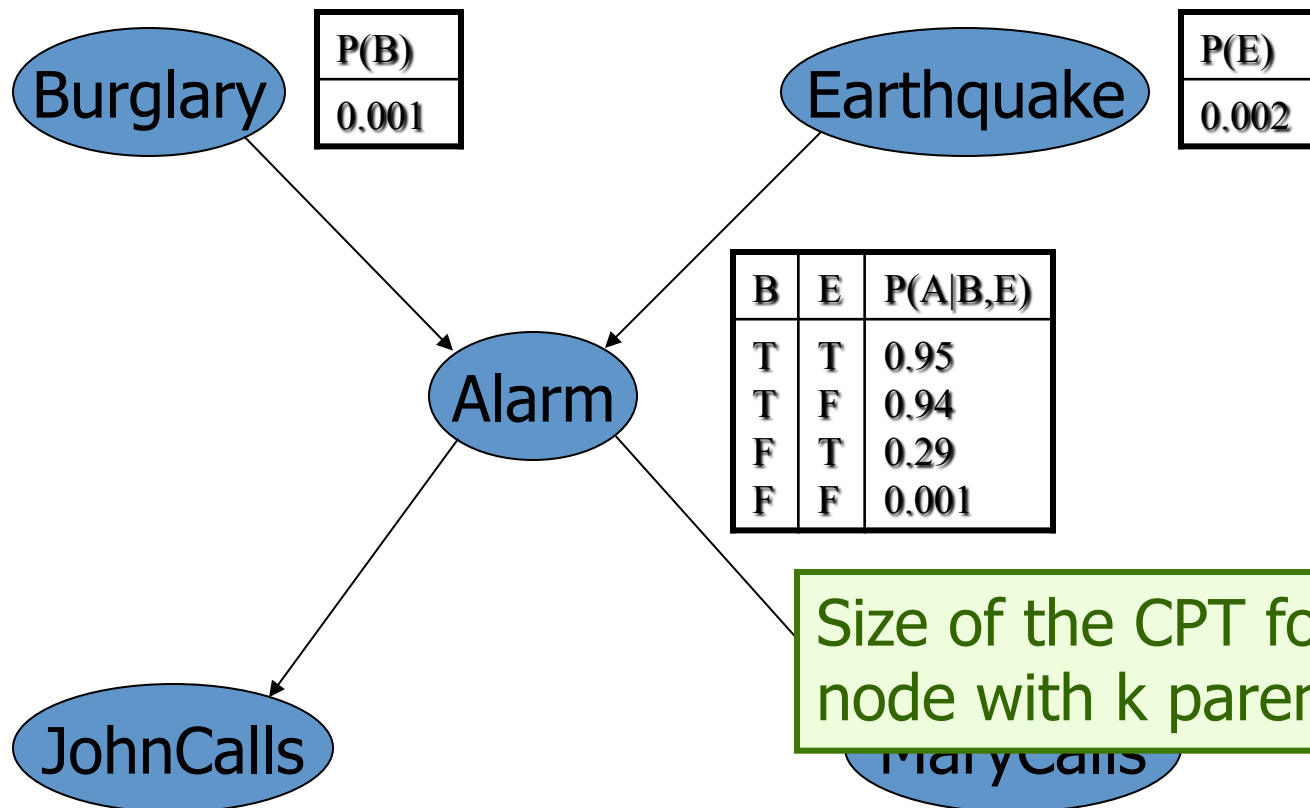- The alarm can cause John to call
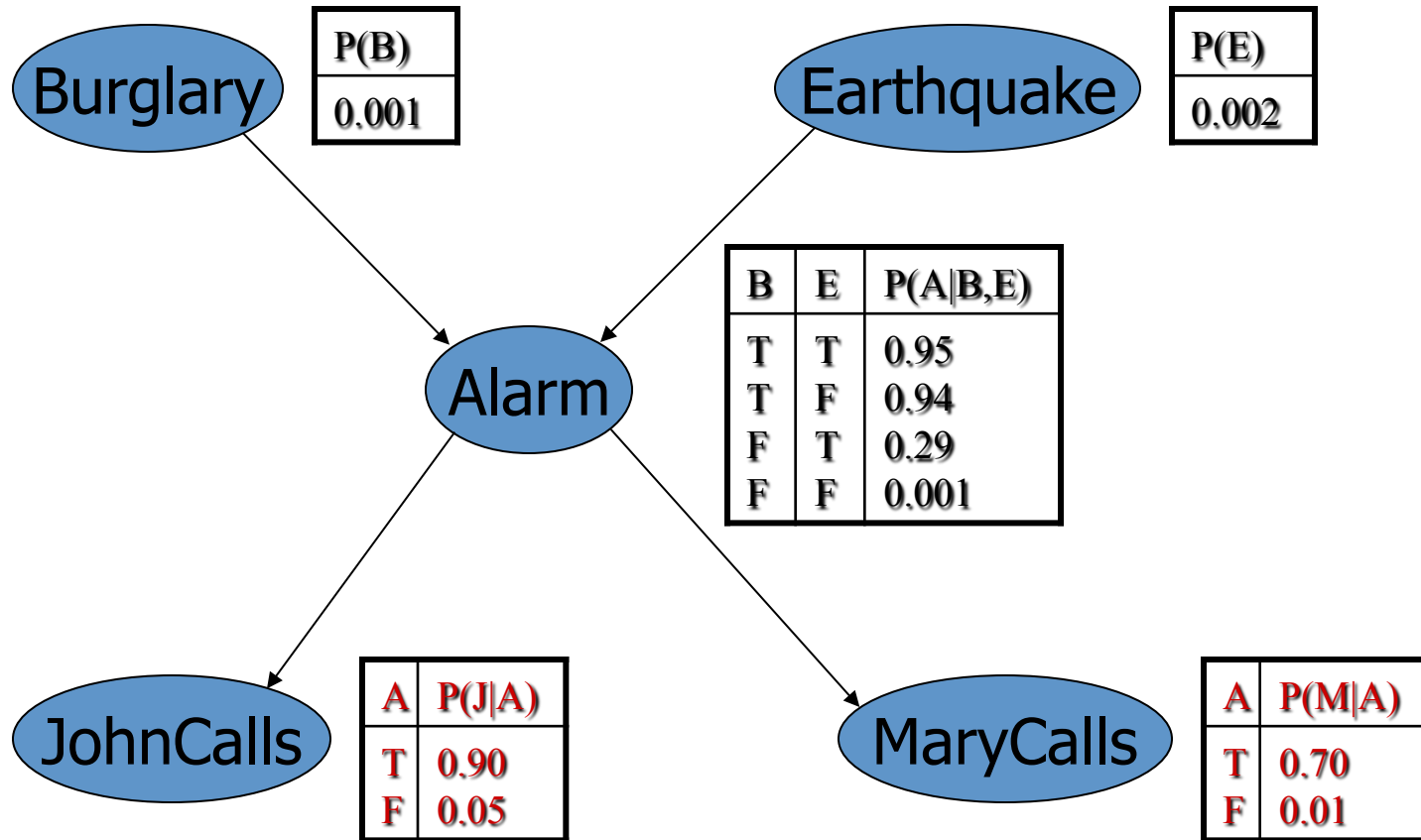
# A Simple Network



Burglary

Earthquake

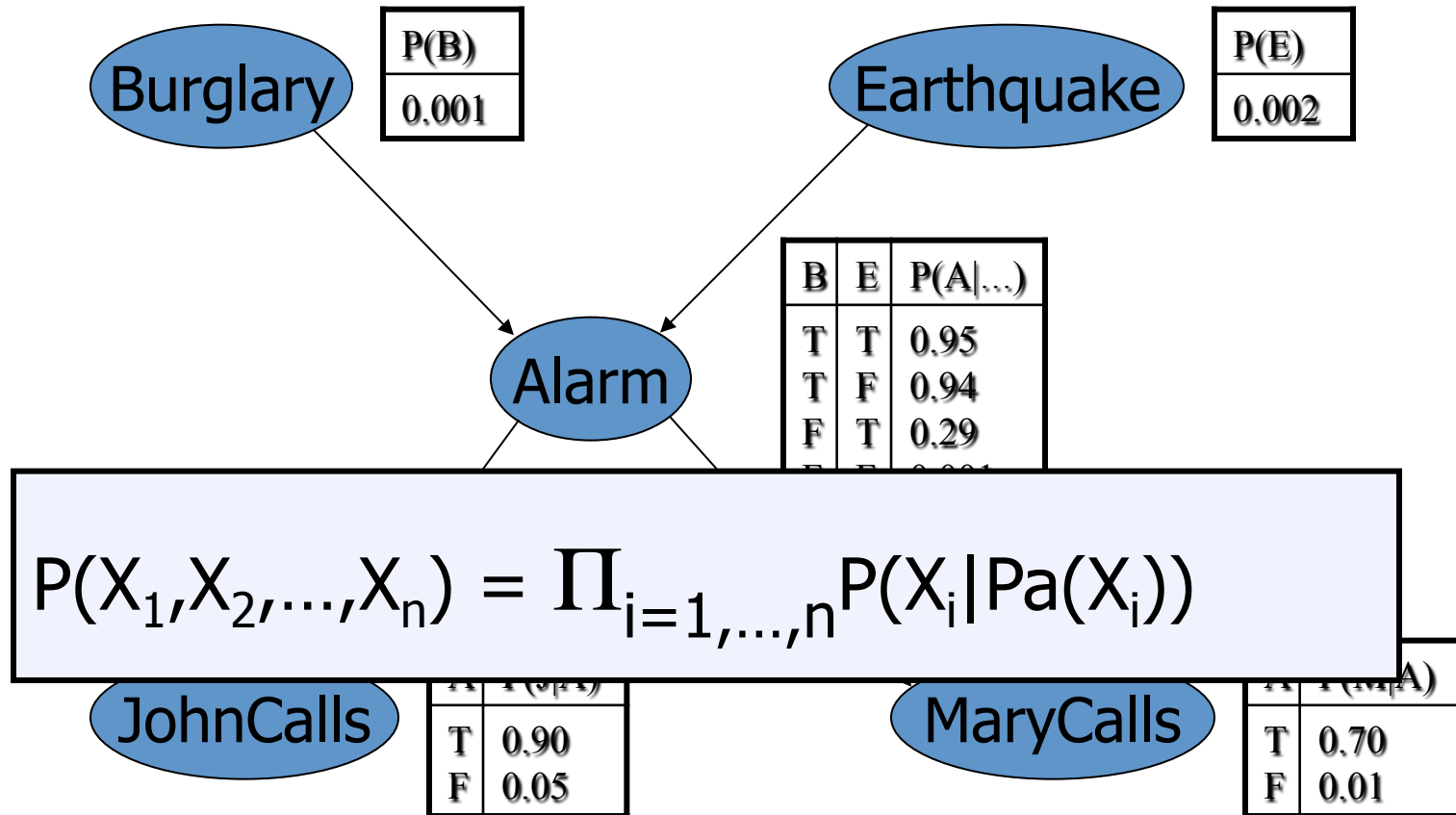Intuitive meaning of arrow from x to y: "x has direct influence on y"

Alarm

Directed acyclic graph (DAG)

JohnCalls

Nodes are random variables

MaryCalls

causes

effects

Colorado State University

# Assigning Probabilities to Roots

# Conditional Probability Tables



| P(B) |
|------|
| 0.001 |

| P(E) |
|------|
| 0.002 |

Burglary          Earthquake

| B | E | P(A\|B,E) |
|---|---|----------|
| T | T | 0.95 |
| T | F | 0.94 |
| F | T | 0.29 |
| F | F | 0.001 |

Alarm

JohnCalls          MaryCalls

Size of the CPT for a node with k parents: ?

# Conditional Probability Tables



**Burglary**

| P(B) |
|------|
| 0.001 |

**Earthquake**

| P(E) |
|------|
| 0.002 |

**Alarm**

| B | E | P(A\|B,E) |
|---|---|---------|
| T | T | 0.95 |
| T | F | 0.94 |
| F | T | 0.29 |
| F | F | 0.001 |

**JohnCalls**

| A | P(J\|A) |
|---|--------|
| T | 0.90 |
| F | 0.05 |

**MaryCalls**

| A | P(M\|A) |
|---|--------|
| T | 0.70 |
| F | 0.01 |

Colorado State University

# What the BN Means



| | P(B) |
|---|---|
| | 0.001 |

| | P(E) |
|---|---|
| | 0.002 |

| B | E | P(A\|...) |
|---|---|---|
| T | T | 0.95 |
| T | F | 0.94 |
| F | T | 0.29 |

$$P(X_1, X_2, ..., X_n) = \Pi_{i=1,...,n} P(X_i | Pa(X_i))$$

JohnCalls

| A | P(J\|A) |
|---|---|
| T | 0.90 |
| F | 0.05 |

MaryCalls

| A | P(M\|A) |
|---|---|
| T | 0.70 |
| F | 0.01 |

Colorado State University

# Calculation of Joint Probability

Burglary

| P(B) |
|------|
| 0.001 |

Earthquake

| P(E) |
|------|
| 0.002 |

| B | E | P(A\|...) |
|---|---|---------|
| T | T | 0.95 |
| T | F | 0.94 |
| F | T | 0.29 |
| F | F | 0.001 |

P(J∧M∧A∧¬B∧¬E)
= P(J|A)P(M|A)P(A|¬B,¬E)P(¬B)P(¬E)
= 0.9 x 0.7 x 0.001 x 0.999 x 0.998
= 0.00062

JohnCalls

| A | P(J\|...) |
|---|---------|
| T | 0.90 |
| F | 0.05 |

MaryCalls

| A | P(M\|...) |
|---|---------|
| T | 0.70 |
| F | 0.01 |

# What the BN Encodes



For example, John does not observe any burglaries directly

- Each of the beliefs JohnCalls and MaryCalls is independent of Burglary and Earthquake given Alarm or ¬Alarm

- The beliefs JohnCalls and MaryCalls are independent given Alarm or ¬Alarm
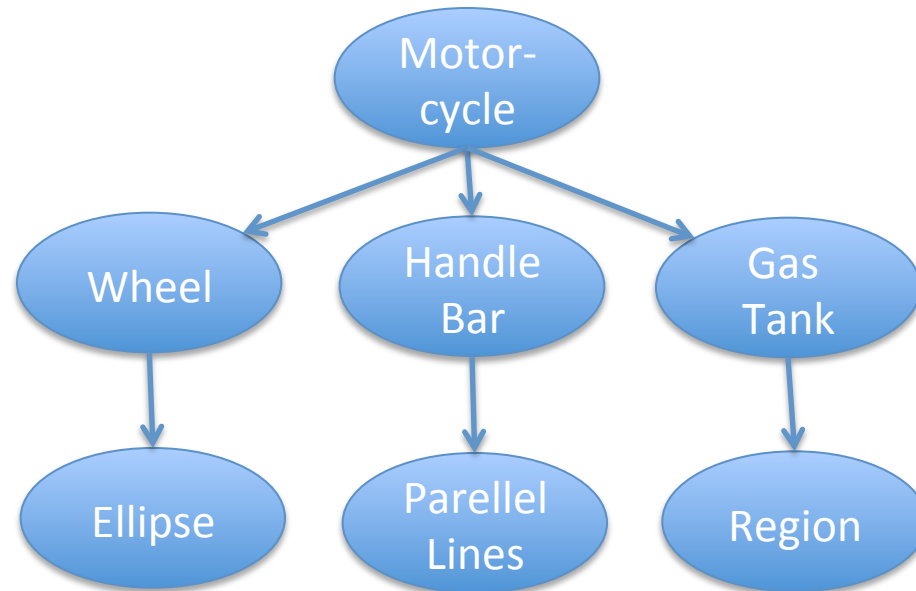
# Independence

- The expression:

$$P(X_1, X_2, \ldots, X_n) = \prod_{i=1,\ldots,n} P(X_i | Pa(X_i))$$

  means that each belief is independent of its predecessors in the BN given its parents

- Said otherwise, the parents of a belief $X_i$ are all the beliefs that "directly influence" $X_i$

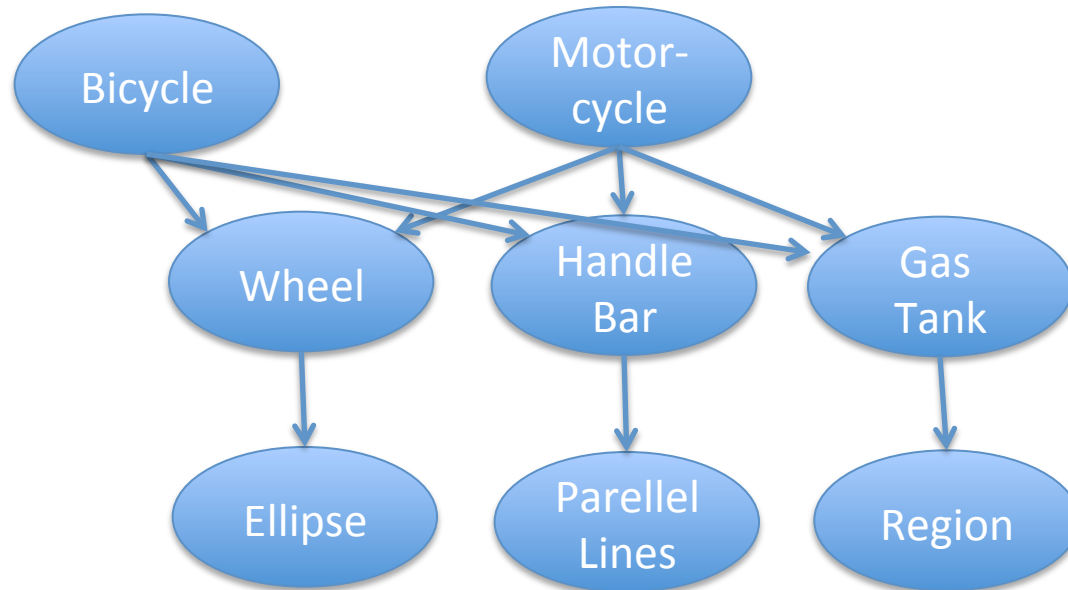- Usually (but not always) the parents of $X_i$ are its causes and $X_i$ is the effect of these causes

E.g., JohnCalls is influenced by Burglary, but not directly. JohnCalls is directly influenced by Alarm
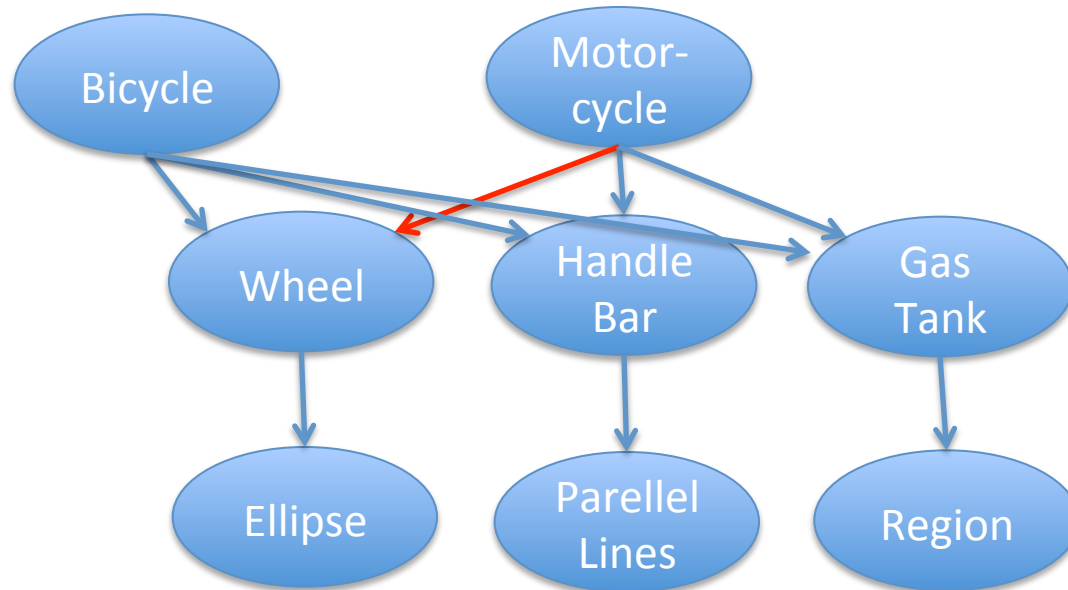
# Computer Vision Example



One purpose for Bayesian nets is to separate object variation information from feature extraction probabilities
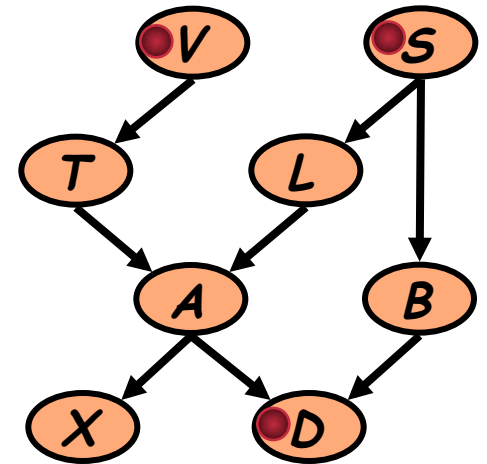
# Computer Vision Example (II)



Another is to exploit shared subpart. Bicycles have wheels and handle bars too, but not gas tanks, so P(Gas Tank | Bicycle) ≈ 0.

# Computer Vision Example (III)



Finally, in the constellation model the probability functions can depend on location as well as existence

Colorado State University

# Solving for Probabilities



- How do we deal with evidence?

- Suppose get evidence $V = t$, $S = f$, $D = t$ and want to compute $P(L|V = t, S = f, D = t)$

$$P(L \mid V = t, S = f, D = t) = \frac{P(L, V = t, S = f, D = t)}{P(V = t, S = f, D = t)}$$

# Variable Elimination

- There is an algorithm for this
  - It's called variable elimination
  - Efficient for trees; inefficient for DAGs
  - Its taught in CS440 (see Russell & Norvig)
- As a user, you need to know:
  - VE can compute the likelihood of any node in a Bayesian net, given any set of evidence
  - But, computing $P(X=x)$ is NP-hard (in general)

# Approaches to inference

- Exact inference
  - Variable elimination
  - Join tree algorithm (also NP)
- Approximate inference
  - Simplify the structure of the network to make exact inference efficient (variational methods, loopy belief propagation)
- Probabilistic methods
  - Stochastic simulation / sampling methods
  - Markov chain Monte Carlo methods