

Machine Learning

11 – Modeling Uncertainty

SS 2018

Gunther Heidemann

Part I: Uncertainty and probability

- Random variables
- Joint distribution
- Inference
- Independence and conditional independence
- Bayes rule

Part II: Bayes networks

Part I: Uncertainty and Probability

Environments may be **uncertain** due to several causes:

- Environment is only partially observable.
- Sensors are unreliable.
- The results of actions are uncertain.
- High complexity.

We deal with uncertainty using **probabilities of propositions**.

Alternative: Model uncertainty using probabilities of rules. such as

$$\begin{aligned} \textit{LawnSprinkler} &\vdash_{0.99} \textit{WetGras}, \\ \textit{WetGras} &\vdash_{0.7} \textit{Rain}. \end{aligned}$$

Probabilities summarize several factors:

- Missing knowledge,
- Incapability to devise complete models of complex domains,
- Chance.

- Modeling uncertainty using **random variables**.
- Types of random variables:
 - **Boolean** random variable:
 - E.g. *Cavity* (Is there a cavity in my tooth?)
 - Values: *<true, false>*
 - **Discrete** random variable:
 - E.g. *Weather* has one of the values *<sun, rain, cloudy, snow>*
 - Values must describe the domain sufficiently and be mutually exclusive.
 - **Continuous** random variable:
 - Values are real numbers.
 - E.g. *Length* $\in [1, 20]$.

- A *proposition* is made by assigning a value to a random variable:
 - *Weather* = *sun*
 - *Length* = 2,4
- Complex propositions are made by using logical operators to connect simple propositions:

$$\textit{Weather} = \textit{sun} \quad \vee \quad \textit{Cavity} = \textit{false}$$

- Notation:
 - Random variables with capital: *Weather*, but values: *sun*.
 - But:

<i>cavity</i>	means	<i>Cavity</i> = <i>true</i> ,
\neg <i>cavity</i>	means	<i>Cavity</i> = <i>false</i> ,
<i>sun</i>	means	<i>Weather</i> = <i>sun</i> .

- **Atomic event:**

A *complete* specification of the state of the domain (the agent may be uncertain about the state).

- Example: Domain is fully described by the boolean variables *Cavity* and *Toothache*.

Then there are 4 atomic events:

Cavity = false \wedge *Toothache = false*

Cavity = false \wedge *Toothache = true*

Cavity = true \wedge *Toothache = false*

Cavity = true \wedge *Toothache = true*

- Atomic events are mutually exclusive and describe the domain completely.

- **A-priori** or **unconditional** probabilities of propositions:
 $P(\text{Cavity} = \text{true}) = 0.1$ or $P(\text{Weather} = \text{sun}) = 0.72$ denote the probability of guesses. The probabilities may change when new information becomes available.
- The **probability distribution** **P** comprises the probabilities of all values:
 $P(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$
 for values $\langle \text{sun}, \text{rain}, \text{cloudy}, \text{snow} \rangle$.
- **P** is **normalized**, i.e., $\text{sum} = 1$.

- **Joint probability distribution** for several random variables comprises all atomic states:

$P(\textit{Weather}, \textit{Cavity})$ is a 4×2 matrix:

<i>Weather</i>	=		<i>sun</i>	<i>rain</i>	<i>cloudy</i>	<i>snow</i>
<i>Cavity</i>	=	true	0.144	0.02	0.016	0.02
<i>Cavity</i>	=	false	0.576	0.08	0.064	0.08

- The joint probability distribution holds the entire knowledge about the domain!

- **Conditional** or **posterior** probability:

E.g. $P(\text{cavity} \mid \text{toothache}) = 0.8$

i.e., the information *toothache* is known (but no more).

- Notation for conditional distributions:

$\mathbf{P}(\text{Cavity} \mid \text{Toothache}) = 2\text{-component vector of 2-comp. vectors}$

- If the additional information *cavity* is known, *Toothache* is irrelevant:

$\mathbf{P}(\text{cavity} \mid \text{Toothache}, \text{cavity}) = \langle 1, 1 \rangle.$

- *sun* is irrelevant given *toothache*:

$P(\text{cavity} \mid \text{toothache}, \text{sun}) = P(\text{cavity} \mid \text{toothache}) = 0.8$

More general:

$\mathbf{P}(\text{Cavity} \mid \text{Toothache}, \text{sun}) = \mathbf{P}(\text{Cavity} \mid \text{Toothache}).$

- Domain knowledge of this kind facilitates finding the joint probability distribution.

- Definition of conditional probability:

$$P(a \mid b) = P(a \wedge b) / P(b) \text{ if } P(b) > 0.$$

- Product rule** is an alternative formulation:

$$P(a \wedge b) = P(a \mid b) P(b) = P(b \mid a) P(a).$$

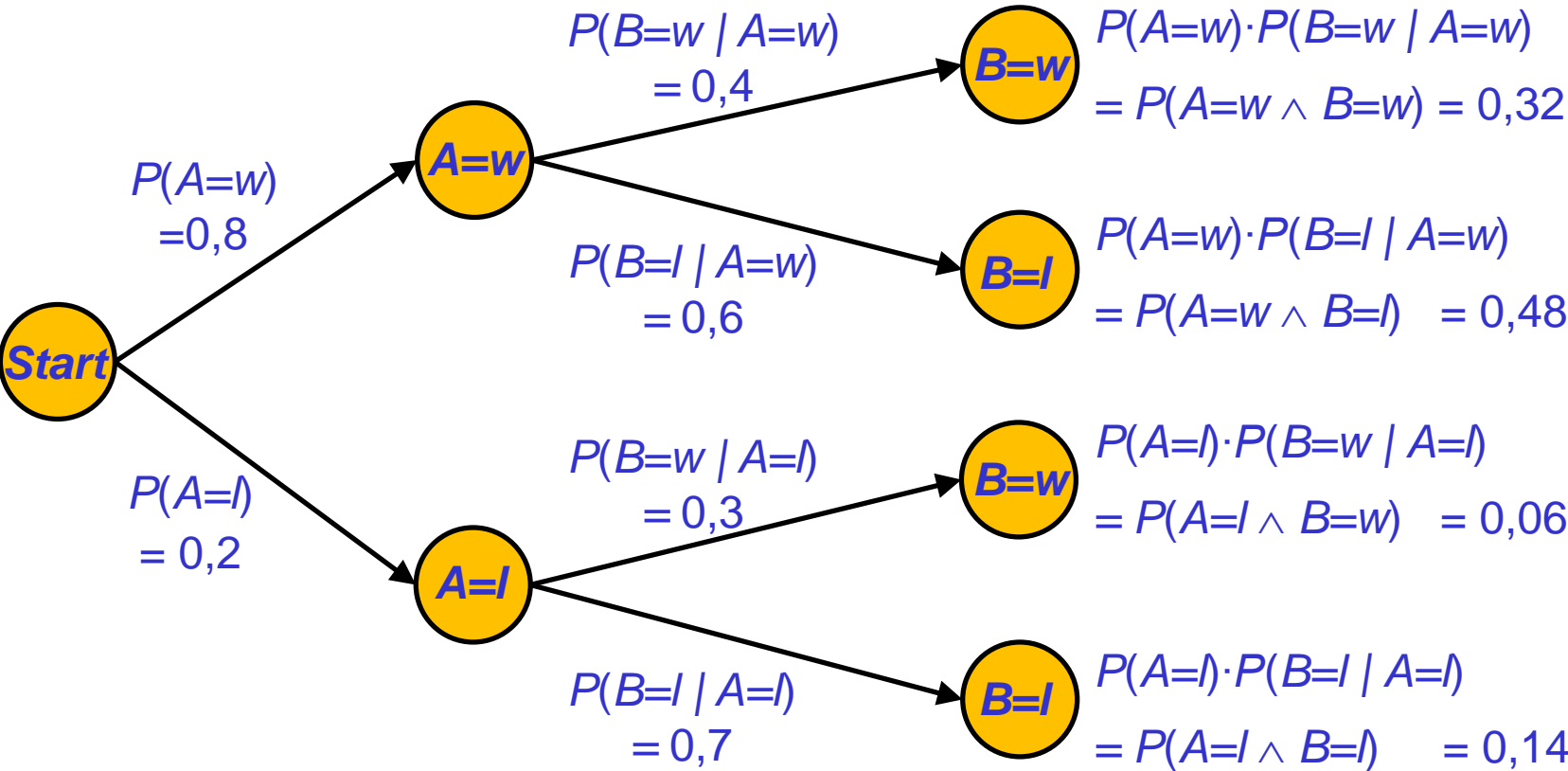
- General version for distributions:

$$\mathbf{P}(\textit{Weather}, \textit{Cavity}) = \mathbf{P}(\textit{Weather} \mid \textit{Cavity}) \mathbf{P}(\textit{Cavity})$$

This means 4×2 separate equations, not matrix multiplication!

Conditional probability

Example: The german national soccer team plays first against Austria (A), winning chance $P(A=w) = 80\%$, then against Brazil (B). Experts say if $A=w$, they will be in a good mood and have a 40% chance to defeat Brazil, otherwise only 30%.



Chain rule (derived by successive application of product rule):

$$\begin{aligned} & \mathbf{P}(X_1, \dots, X_n) \\ &= \mathbf{P}(X_n \mid X_1, \dots, X_{n-1}) \mathbf{P}(X_1, \dots, X_{n-1}) \\ &= \mathbf{P}(X_n \mid X_1, \dots, X_{n-1}) \mathbf{P}(X_{n-1} \mid X_1, \dots, X_{n-2}) \mathbf{P}(X_1, \dots, X_{n-2}) \\ &= \dots \\ &= \prod_{i=2}^n \mathbf{P}(X_i \mid X_1, \dots, X_{i-1}) \mathbf{P}(X_1). \end{aligned}$$

Joint probability distribution (*catch* = dentist has found cavity):

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

[RN]

- Probability of a proposition ϕ is the sum of the probabilities of the corresponding atomic events:

$$P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega).$$

- $P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$
- $P(\text{toothache} \vee \text{cavity}) =$
 $0.108 + 0.012 + 0.016 + 0.064 + 0.072 + 0.008 = 0.28.$

Joint probability distribution (*catch* = dentist has found cavity):

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

[RN]

Conditional probabilities:

$$P(\neg \text{cavity} \mid \text{toothache}) =$$

$$P(\neg \text{cavity} \wedge \text{toothache}) / P(\text{toothache})^* =$$

$$(0.016 + 0.064) / (0.108 + 0.012 + 0.016 + 0.064) = 0.4$$

* Since *toothache* is known, the right side of the table must be normalized to 1.

Joint probability distribution (*catch = dentist has found cavity*):

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

[RN]

Denominator is a **normalization constant**:

$$\alpha = 1 / P(\text{toothache})$$

$P(\text{Cavity} \mid \text{toothache})$

$$= \alpha P(\text{Cavity}, \text{toothache})$$

$$= \alpha [P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \neg \text{catch})]$$

$$= \alpha [<0.108, 0.016> + <0.012, 0.064>]$$

$$= \alpha <0.12, 0.08> = <0.6, 0.4>$$

Idea: Infer distribution of *query variables* (*Cavity*) depending on the *evidence variables* (*toothache*) and sum up over unobserved *hidden variables* (*Catch*).

In general:

For a set X of random variables we want to know

- the joint posterior distributions of the *query variables* Y
- for given values e of the *evidence variables* E .

The *hidden variables* are $H = X \setminus \{Y \cup E\}$,

and can be removed by “summing out”:

$$P(Y \mid E = e) = \alpha P(Y, E = e) = \alpha \sum_h P(Y, E=e, H=h)$$

Problem of inference by enumeration:

For n variables with a maximum of d values, the joint probability distribution table comprises $O(d^n)$ values:

- How to find these numbers?
- Memory requirements are $O(d^n)$.
- Time complexity is $O(d^n)$.

Domain knowledge about independence of variables may simplify the joint probability distribution significantly !

A and B are **independent** if and only if

$$P(A \mid B) = P(A) \quad \text{or} \quad P(B \mid A) = P(B).$$

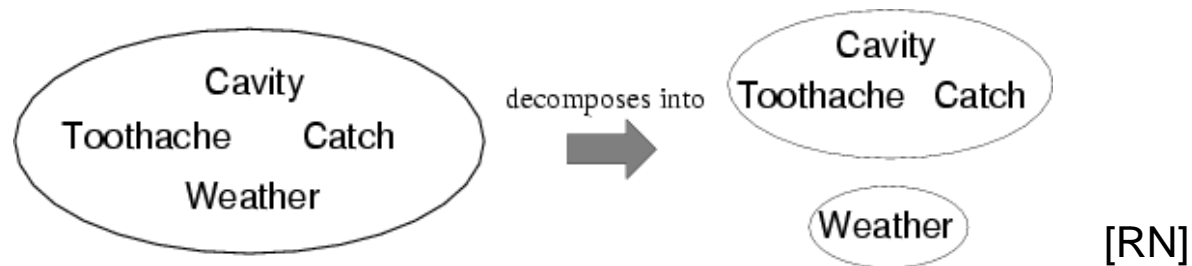
With this, the product rule leads to

$$P(A, B) = P(A \mid B) P(B) = P(A) P(B).$$

Example:

$$\begin{aligned}
 & P(\textit{Toothache}, \textit{Cavity}, \textit{Catch}, \textit{Weather}) \\
 = & P(\textit{Weather} \mid \textit{Toothache}, \textit{Cavity}, \textit{Catch}) P(\textit{Toothache}, \textit{Cavity}, \textit{Catch}) \\
 = & P(\textit{Weather}) P(\textit{Toothache}, \textit{Cavity}, \textit{Catch}),
 \end{aligned}$$

since *Weather* is independent of *Toothache*, *Cavity*, *Catch*.



So the $2 \times 2 \times 2 \times 4 - 1 = 31$ independent (sum = 1 !) numbers of the joint probability distribution are reduced to $4 + 2 \times 2 \times 2 - 1 = 11$.

But: Perfect independence is rare (toothache might be influenced by the weather).

Consider $P(\textit{Toothache}, \textit{Cavity}, \textit{Catch})$, which has $2^3 - 1 = 7$ independent probabilities.

Mind \textit{Catch} is **not independent** of $\textit{Toothache}$:

In general,

$$P(\textit{Catch} \mid \textit{Toothache}) \neq P(\textit{Catch}) !$$

Rather, $P(\textit{Catch})$ **does** depend on the value of $\textit{Toothache}$ - it's far more likely finding a cavity provided there is $\textit{toothache}$.

But: This holds only as long as the value of \textit{Cavity} is unknown.

Assume that

1. for $Cavity = true$ the probability that the dentist finds the cavity does not depend on whether there is toothache or not

$$P(Catch \mid Toothache, cavity) = P(Catch \mid cavity);$$

2. for $Cavity = false$ the probability for a catch is independent of $Toothache$ likewise:

$$P(Catch \mid Toothache, \neg cavity) = P(Catch \mid \neg cavity).$$

Summarizing 1 and 2, $Catch$ is **conditionally independent** of $Toothache$ **given the value of** $Cavity$:

$$P(Catch \mid Toothache, Cavity) = P(Catch \mid Cavity).$$

Likewise, $Toothache$ is conditionally independent of $Catch$ given $Cavity$:

$$P(Toothache \mid Catch, Cavity) = P(Toothache \mid Cavity)$$

(e.g., the tooth hurts whether the dentist finds the cavity or not).

From the conditional independence of *Catch* and *Toothache* for given *Cavity*

$$P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$$

$$P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity})$$

we get the factorization

$$\begin{aligned} &P(\text{Toothache}, \text{Catch} \mid \text{Cavity}) \\ &= P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) \cdot P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) \\ &= P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity}). \end{aligned}$$

Full joint probability distribution derived using chain rule:

$$\begin{aligned}
 &P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) \\
 &= P(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity}) \cdot P(\textit{Catch}, \textit{Cavity}) \\
 &= P(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity}) \cdot P(\textit{Catch} \mid \textit{Cavity}) P(\textit{Cavity}) \\
 &= P(\textit{Toothache} \mid \textit{Cavity}) \quad \cdot \quad P(\textit{Catch} \mid \textit{Cavity}) P(\textit{Cavity}),
 \end{aligned}$$

i.e., $2 + 2 + 1 = 5$ independent numbers.

- In many cases, the memory required to represent a joint probability distribution of n variables can be reduced from “exponential in n ” to “linear in n ” using knowledge about conditional independence.
- Thus, conditional independence is an important and simple method of representing knowledge.

- Product rule: $P(A, B) = P(A | B) \cdot P(B) = P(B | A) \cdot P(A)$
 \Rightarrow Bayes rule: $P(A | B) = P(B | A) \cdot P(A) / P(B)$
- The same for distributions:
$$P(Y | X) = P(X | Y) P(Y) / P(X) = \alpha P(X | Y) P(Y).$$
- Useful for assessing *diagnostic* probability from *causal* probability :
 - $P(\text{Cause} | \text{Effect}) =$
$$P(\text{Effect} | \text{Cause}) P(\text{Cause}) / P(\text{Effect})$$
 - Example: Let *M* be meningitis, *S* stiff neck
$$P(m | s) = P(s | m) P(m) / P(s) = 0.8 \times 0.0001 / 0.1 = 0.0008$$

The posterior probability of meningitis is small even for a stiff neck, because the a-priori probability for meningitis is small while the a-priori probability for a stiff neck is much larger.

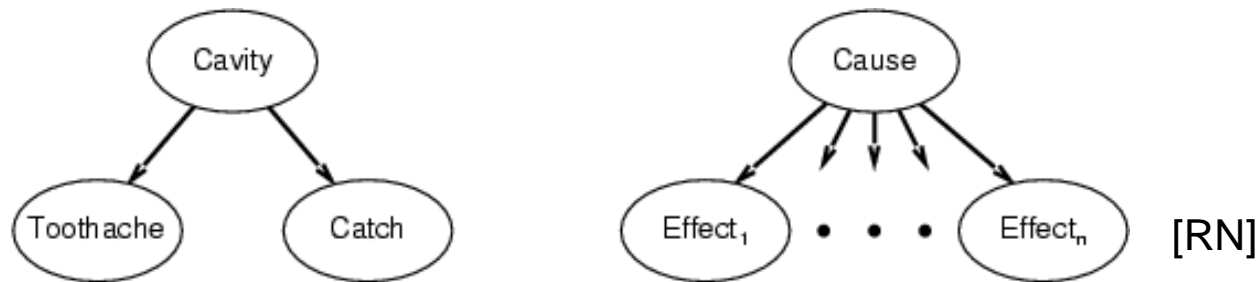
Naive Bayes model:

Let denote C the cause, E_i its effects.

$$\begin{aligned}
 P(C, E_1, \dots, E_n) &= P(E_1 | C, E_2 \dots E_n) P(C, E_2 \dots E_n) \\
 &= P(E_1 | C, E_2 \dots E_n) P(E_2 | C, E_3 \dots E_n) P(C, E_3 \dots E_n) \\
 \text{etc.} \quad &= \prod_{i=1}^{n-1} P(E_i | C, E_{i+1} \dots E_n) P(E_n | C) P(C) \\
 &= \prod_{i=1}^n P(E_i | C) P(C)
 \end{aligned}$$

since effect E_i is conditionally independent of the other effects $E_j, j \neq i$, for a given value of cause C .

The number of parameters is linear in n .



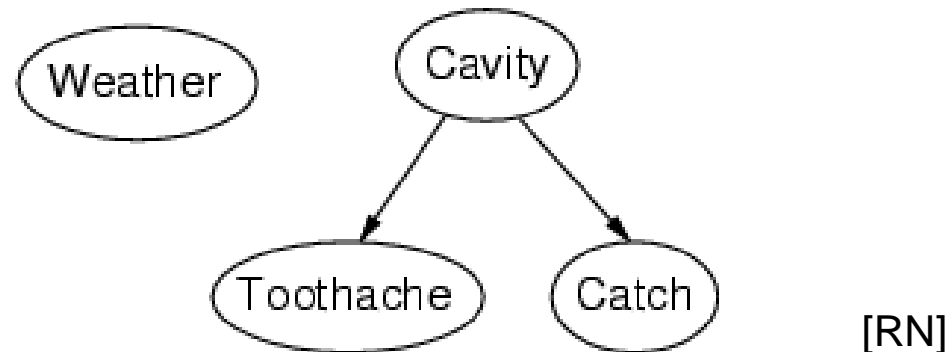
Part II: Bayes networks

- **Bayes networks** are a form of graphical notation for propositions on conditional probabilities and thus a suitable way to express joint probability distributions.
- Syntax:
 - A set of **nodes**, one for each random variable.
 - A directed **acyclic graph**
 - Edge from A to B means: „ A influences B “.
 - A is a **parent node** of B .
 - A conditional probability distribution for each node depending on its parent nodes:

$$P(X_i | Parents(X_i))$$

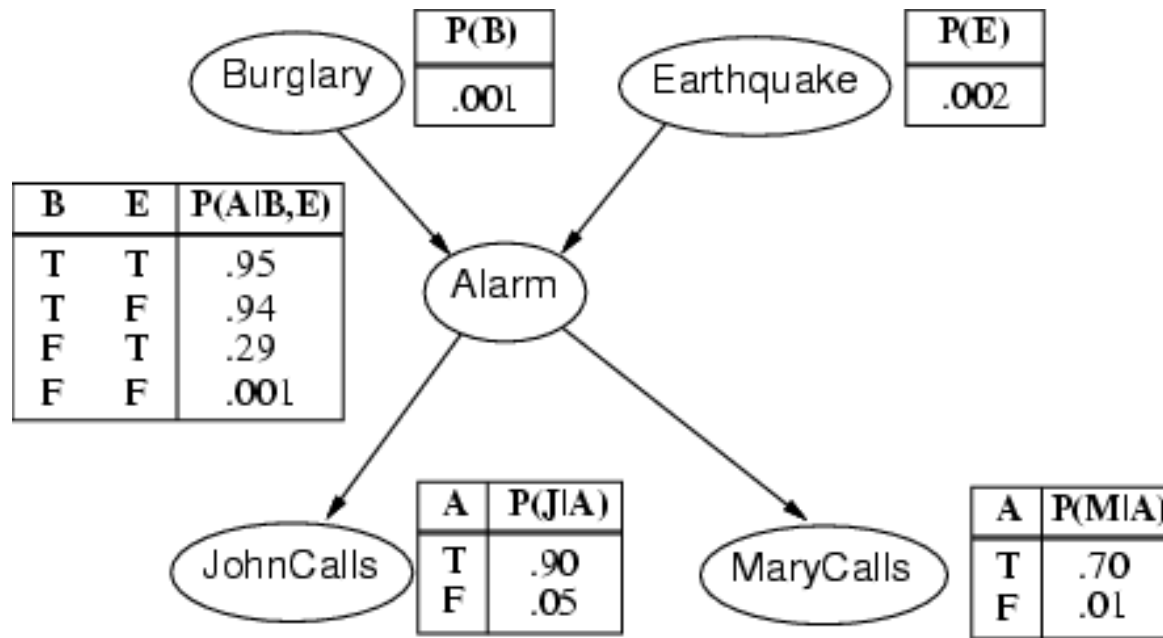
- In the simplest case, the conditional probability distribution is represented as a **conditional probability table (CPT)**, which specifies the distribution over X_i for each combination of parent values.

The topology of the net encodes conditional independence assertions:



- *Weather* is independent of the other variables.
- *Toothache* und *Catch* are conditionally independent given the value of *Cavity*.

- I'm not at home. My neighbor John calls because my alarm is ringing but neighbor Mary does not call. The alarm should indicate burglars, but sometimes it is set off by minor earthquakes. Is there a burglary?
- Variables: *Burglary, Earthquake, Alarm, JohnCalls, MaryCalls*
- The joint probability distribution of the five variables comprises $2^5 - 1 = 31$ independent numbers (without further knowledge about independencies).
- **Network topology reflects causal knowledge:**
 - A burglary can set off the alarm.
 - An earthquake can set off the alarm.
 - The alarm can cause Mary to call.
 - The alarm can cause John to call.



[RN]

Independencies: E.g., *MaryCalls* is **not independent** on *JohnCalls* (because the alarm makes Mary's calling likely, but the same applies to John).

But *MaryCalls* is **conditionally independent** on *JohnCalls* given *Alarm*.

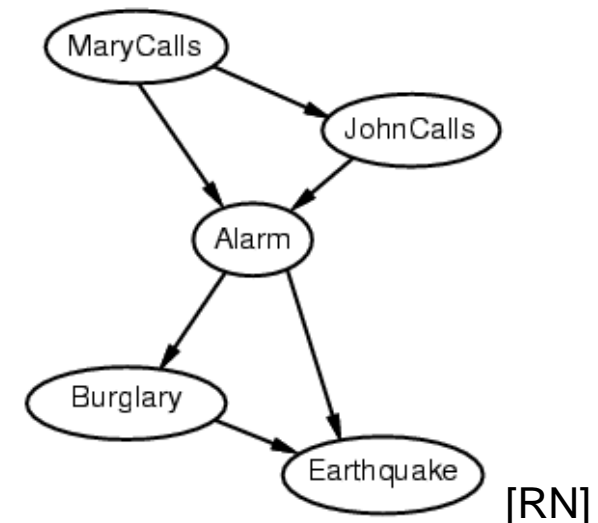
Knowledge about conditional independencies reduces joint probability distribution to $1 + 1 + 4 + 2 + 2 = 10$ independent numbers.

The joint probability distribution is the product of the local conditional probability distributions:

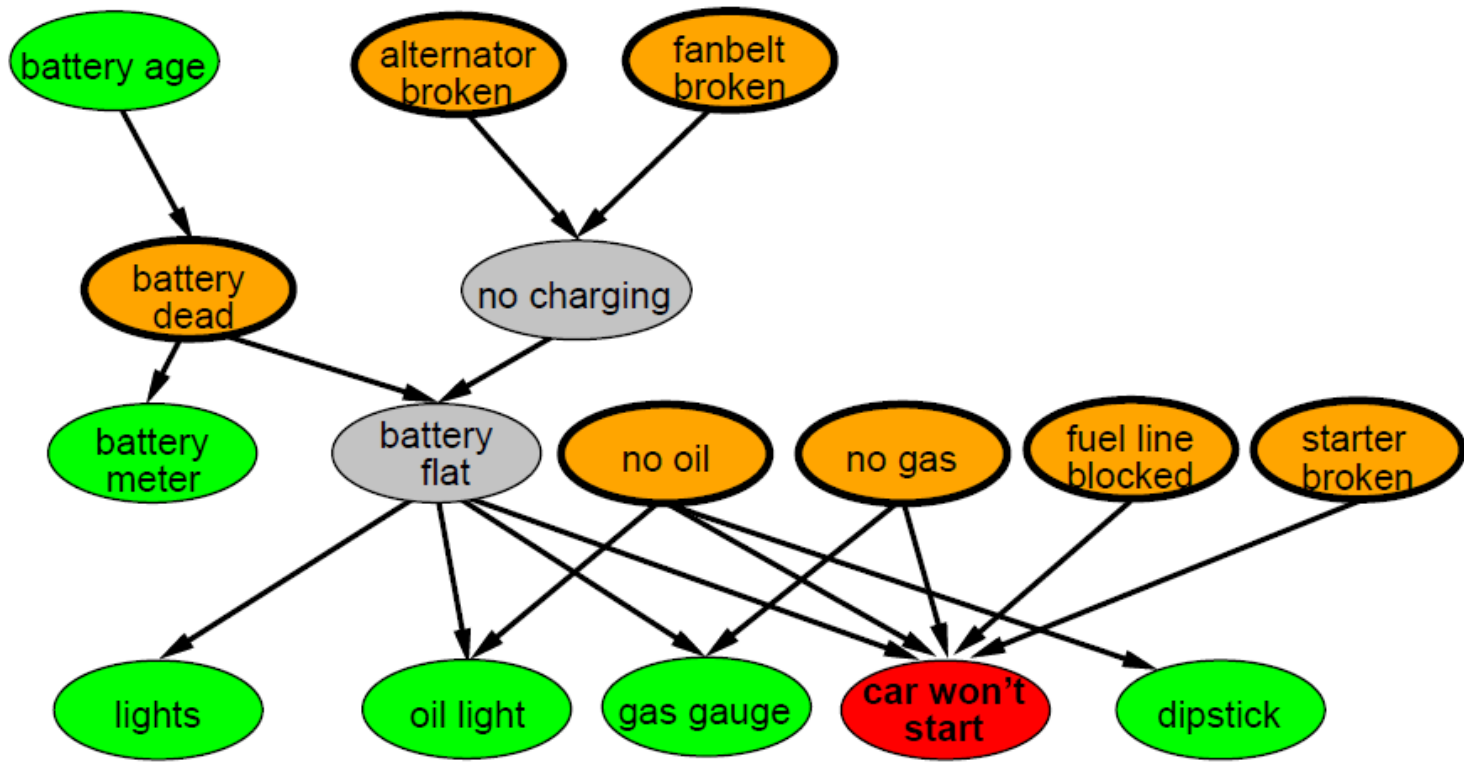
$$\mathbf{P}(X_1, \dots, X_n) = \prod_{i=1..n} \mathbf{P}(X_i | \text{Parents}(X_i)).$$

Note:

- It is convenient to construct the edges in the causal direction, but the opposite (diagnostic direction) is possible as well.
- However, dependencies change and so may
- the number of independent parameters, and
- the network becomes more difficult to interpret.



Example



[RN]

- Uncertainties can be represented by assigning probabilities to propositions (but that's not the only way).
- The joint probability distribution assigns a probability to each atomic event and thus holds the complete domain knowledge.
- Inference is possible by summing up probabilities of atomic events.
- Independence reduces the complexity of the joint probability distribution but is rarely perfect in reality.
- Conditional independence is more feasible.
- Bayes nets are a convenient representation for the dependence / conditional independence of random variables.
- Causal direction of edges is easier to interpret.
- Bayes nets are easier to design for domain experts than conditional probabilities.

- [M] Online material available at www.cs.cmu.edu/~tom/mlbook.html for the textbook: Tom M. Mitchell: *Machine Learning*, McGraw-Hill
- [RN] Stuart Russell, Peter Norvig: *Artificial Intelligence*, Pearson
- [H] Gunther Heidemann, 2012.