# Reasoning Under Uncertainty: Bayesian networks intro

CPSC 322 – Uncertainty 4

Textbook §6.3 – 6.3.1

March 23, 2011

# Lecture Overview

Recap: marginal and conditional independence

- Bayesian Networks Introduction

- Hidden Markov Models
  - Rainbow Robot Example

# Marginal Independence

**Definition (Marginal independence)**
Random variable X is (marginally) independent of random variable Y, written X ⫫ Y, if for all x ∈ dom(X), $y_j$ ∈ dom(Y) and $y_k$ ∈ dom(Y), the following equation holds:

$$P(X = x | Y = y_j)$$
$$= P(X = x | Y = y_k)$$
$$= P(X = x)$$

- Intuitively: if X ⫫ Y, then
  - learning that Y=y does not change your belief in X
  - and this is true for all values y that Y could take

- For example, weather is marginally independent from the result of a coin toss

# Marginal Independence

**Definition (Marginal independence)**
Random variable X is (marginally) independent of random variable Y, written X ⫫ Y, if for all x ∈ dom(X), $y_j$ ∈ dom(Y) and $y_k$ ∈ dom(Y), the following equation holds:

$$P(X = x | Y = y_j)$$
$$= P(X = x | Y = y_k)$$
$$= P(X = x)$$

- Recall the product rule:
  - $P(X = x \land Y = y) = P(X = x | Y = y) \times P(Y = y)$

- If X ⫫ Y, we have:
  - $P(X = x \land Y = y) = P(X = x) \times P(Y = y)$
  - In distribution form: $P(X, Y) = P(X) \times P(Y)$

- If $X_i$ ⫫ $X_j$ for all i, j:  $P(X_1, ..., Xn) = \prod_{i=1}^{n} P(Xi)$

# Conditional Independence

**Definition (Conditional independence)**

Random variable X is (conditionally) independent of random variable Y given random variable Z, written $X \perp\!\!\!\perp Y \mid Z$ if, for all $x \in dom(X)$, $y_j \in dom(Y)$, $y_k \in dom(Y)$ and $z \in dom(Z)$ the following equation holds:

$$P(X = x \mid Y = y_j, Z = z)$$
$$= P(X = x \mid Y = y_k, Z = z)$$
$$= P(X = x \mid Z = z)$$

- Intuitively: if $X \perp\!\!\!\perp Y \mid Z$, then
  - learning that Y=y does not change your belief in X when we already know Z=z
  - and this is true for all values y that Y could take and all values z that Z could take

- For example,
  ExamGrade $\perp\!\!\!\perp$ AssignmentGrade | UnderstoodMaterial

# Conditional Independence

**Definition (Conditional independence)**

Random variable X is (conditionally) independent of random variable Y given random variable Z, written X ⊥⊥ Y | Z if, for all x ∈ dom(X), $y_j$ ∈ dom(Y), $y_k$ ∈ dom(Y) and z ∈ dom(Z) the following equation holds:

$$P(X = x | Y = yj, Z = z)$$
$$= P(X = x | Y = yk, Z = z)$$
$$= P(X = x | Z = z)$$

- Definition of X ⊥⊥ Y | Z in distribution form: $P(X|Y, Z) = P(X|Z)$

- Product rule still holds when every term is conditioned on Z=z:
  - $P(X = x \land Y = y | Z = z) = P(X = x | Y = y, Z = z) \times P(Y = y | Z = z)$

- Thus, if X ⊥⊥ Y | Z :
  - $P(X = x \land Y = y | Z = z) = P(X = x | Z = z) \times P(Y = y | Z = z)$
  - In distribution form: $P(X, Y | Z) = P(X|Z) \times P(Y|Z)$
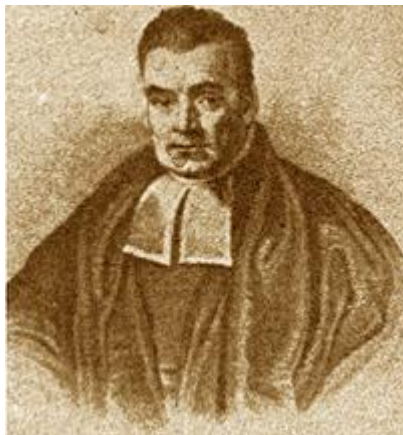
# Lecture Overview

- Recap: marginal and conditional independence
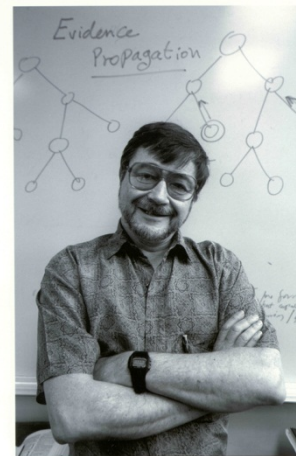
Bayesian Networks Introduction

- Hidden Markov Models
  - Rainbow Robot Example

# Bayesian Network Motivation

- We want a representation and reasoning system that is based on conditional (and marginal) independence
  - Compact yet expressive representation
  - Efficient reasoning procedures
- Bayesian Networks are such a representation
  - Named after Thomas Bayes (ca. 1702 –1761)
  - Term coined in 1985 by Judea Pearl (1936 –  )
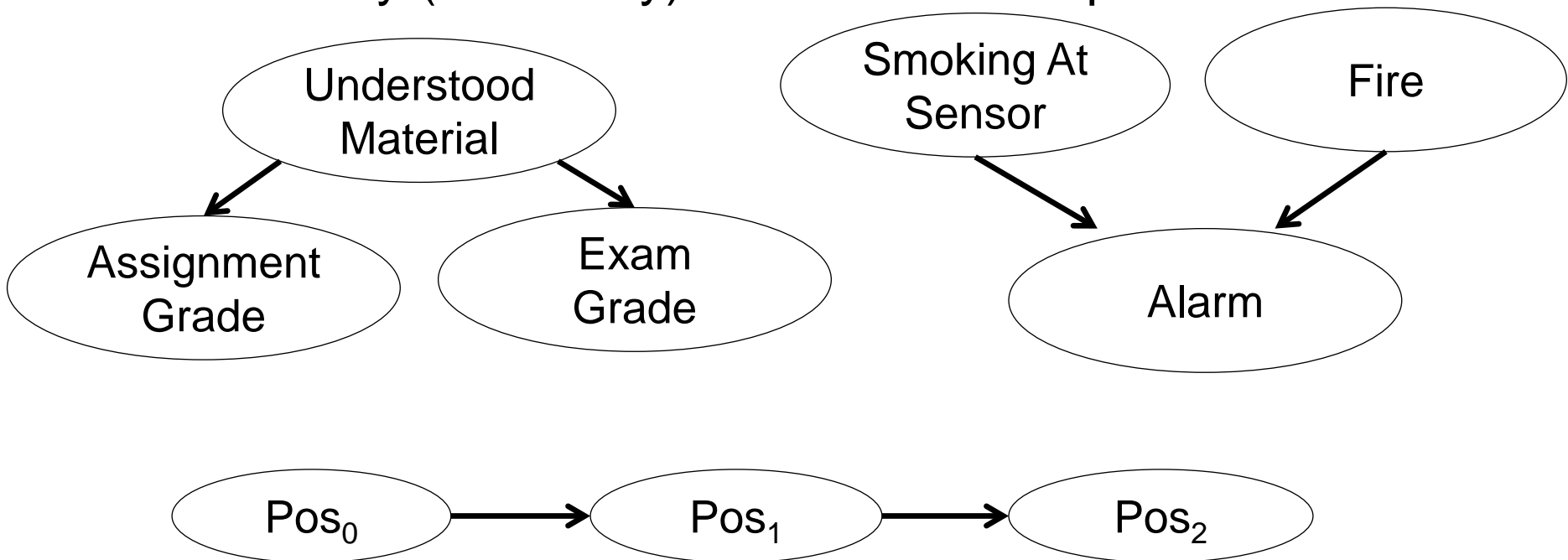  - Their invention changed the focus on AI from logic to probability!

Thomas Bayes                    Judea Pearl

# Bayesian Networks: Intuition

- A graphical representation for a joint probability distribution
    - Nodes are random variables
    - Directed edges between nodes reflect dependence

- We already (informally) saw some examples:

Understood Material → Assignment Grade

Understood Material → Exam Grade

Smoking At Sensor → Alarm

Fire → Alarm

$Pos_0$ → $Pos_1$ → $Pos_2$

# Bayesian Networks: Definition

**Definition (Bayesian Network)**

A Bayesian network consists of

- A directed acyclic graph $(V, E)$ whose nodes are labeled with random variables
- A domain for each random variable
- A conditional probability distribution for each variable $V$
  - Specifies $P(V|Parents(V))$
  - $Parents(V)$ is the set of variables $V'$ with $(V', V) \in E$
    - For nodes $V$ without predecessors, $Parents(V) = \{\}$

- The parents of variable $V$ are those $V$ directly depends on

- A Bayesian network is a compact representation of the JPD: $P(X_1, \ldots, X_n) = \prod_{i=1}^{n} P(X_i|Pa(X_i))$

- Other names for Bayesian networks:
  - Bayes nets, Belief networks, Bayesian Belief networks
  - Common abbreviation: BN

# Bayesian Networks: Definition

**Definition (Bayesian Network)**

A Bayesian network consists of
- A directed acyclic graph $(V, E)$ whose nodes are labeled with random variables
- A domain for each random variable
- A conditional probability distribution for each variable $V$
  - Specifies $P(V|Parents(V))$
  - $Parents(V)$ is the set of variables $V'$ with $(V', V) \in E$
    - For nodes $V$ without predecessors, $Parents(V) = \{\}$

- Discrete Bayesian networks:
  - Domain of each variable is finite
  - Conditional probability distribution is a conditional probability table
  - We will assume this discrete case
    - But everything we say about independence (marginal & conditional) carries over to the continuous case

11

# Example for BN construction: Fire Diagnosis

Bayesian networks are a compact representation of the joint probability distribution (over all variables in the network)
Encoding the joint over $\mathcal{X}$ = {$X_1$, …, $X_n$} as a Bayesian network:

1. Totally order the variables of interest: $X_1$, …, $X_n$

2. Use chain rule with that ordering: $P(X_1, …, X_n) = \prod_{i=1}^{n} P(X_i|X_{i-1},...,X_1)$

3. For every variable $X_i$, find the smallest set of parents
   $Pa(X_i) \subseteq$ {$X_1$, …, $X_{i-1}$} such that $X_i \perp\!\!\!\perp$ {$X_1$, …, $X_{i-1}$} \ $Pa(X_i)$ | $Pa(X_i)$
   - $X_i$ is conditionally independent from its other ancestors given its parents

4. Then we can rewrite $P(X_1, …, X_n) = \prod_{i=1}^{n} P( X_i|Pa(X_i) )$
   - This is a compact representation of the joint probability distribution

5. Construct the BN
   - Nodes are variables
   - Directed edges from all variables in $Pa(X_i)$ to $X_i$
   - Conditional probability table for each variable $X_i$ : $P(X_i | Pa(X_i) )$

12

# Example for BN construction: Fire Diagnosis

You want to diagnose whether there is a fire in a building

- You receive a noisy report about whether everyone is leaving the building

- If everyone is leaving, this may have been caused by a fire alarm

- If there is a fire alarm, it may have been caused by a fire or by tampering

- If there is a fire, there may be smoke

# Example for BN construction: Fire Diagnosis

First you choose the variables. In this case, all are Boolean:

- Tampering is true when the alarm has been tampered with

- Fire is true when there is a fire

- Alarm is true when there is an alarm

- Smoke is true when there is smoke

- Leaving is true if there are lots of people leaving the building

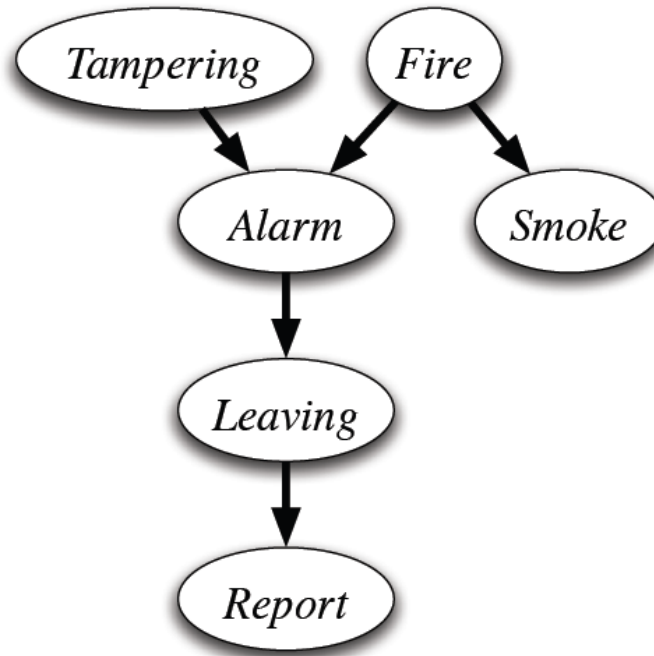- Report is true if the sensor reports that lots of people are leaving the building


- Let's construct the Bayesian network for this (whiteboard)
  - First, you choose a total ordering of the variables, let's say: Fire; Tampering; Alarm; Smoke; Leaving; Report.

# Example for BN construction: Fire Diagnosis

- Using the total ordering of variables:
  - Let's say Fire; Tampering; Alarm; Smoke; Leaving; Report.

- Now choose the parents for each variable by evaluating conditional independencies
  - Fire is the first variable in the ordering. It does not have parents.
  - Tampering independent of fire (learning that one is true would not change your beliefs about the probability of the other)
  - Alarm depends on both Fire and Tampering: it could be caused by either or both
  - Smoke is caused by Fire, and so is independent of Tampering and Alarm given whether there is a Fire
  - Leaving is caused by Alarm, and thus is independent of the other variables given Alarm
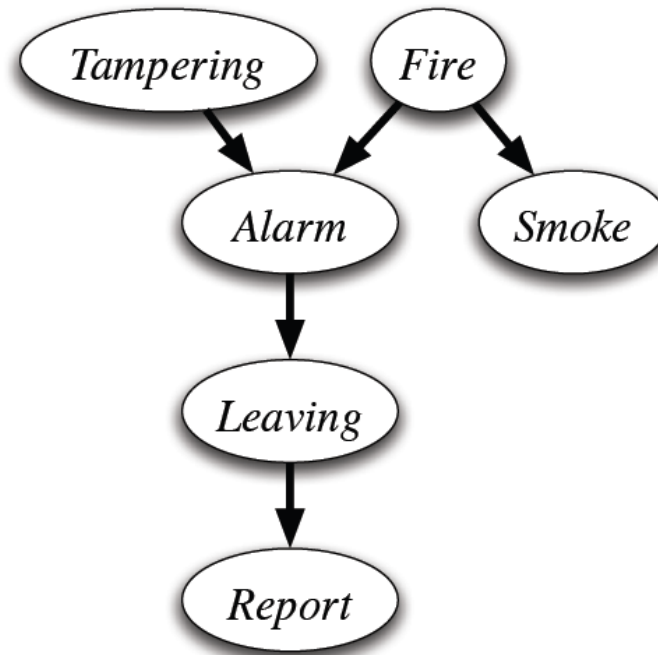  - Report is caused by Leaving, and thus is independent of the other variables given Leaving

# Example for BN construction: Fire Diagnosis

- This results in the following Bayesian network



- P(Tampering, Fire, Alarm, Smoke, Leaving, Report)
  = P(Tampering) × P(Fire) × P(Alarm|Tampering,Fire)
    × P(Smoke|Fire) × P(Leaving|Alarm) × P(Report|Leaving)
- Of course, we're not done until we also come up with conditional probability tables for each node in the graph
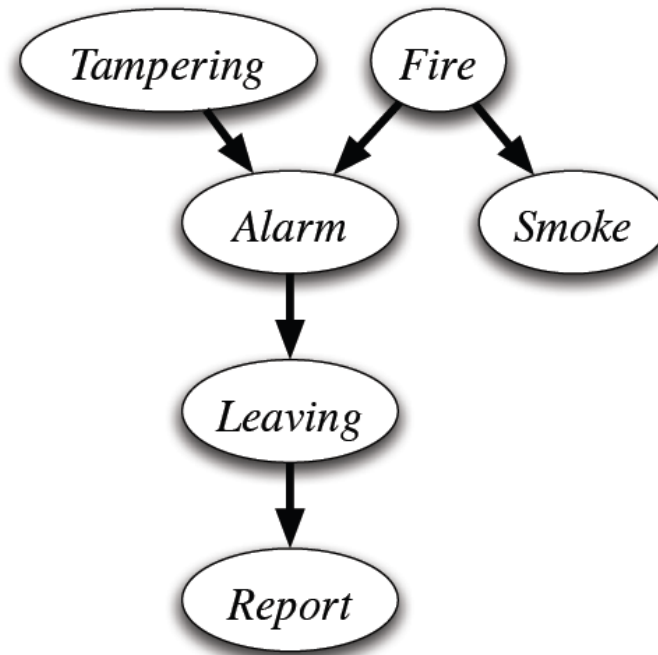
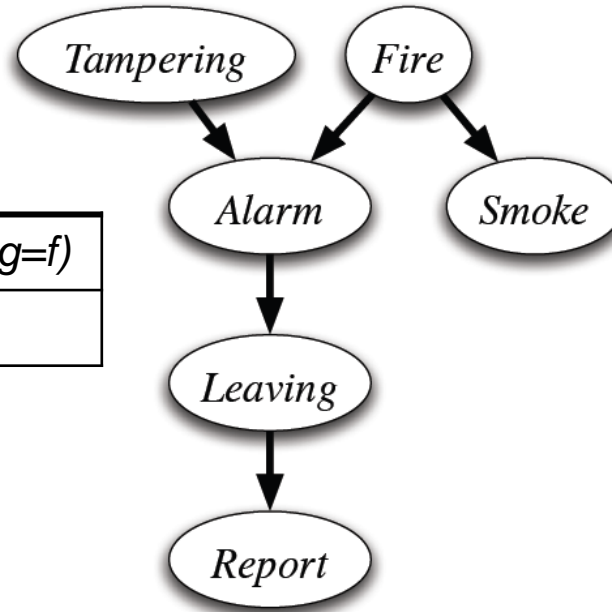# Example for BN construction: Fire Diagnosis



- All variables are Boolean

- How many probabilities do we need to specify for this Bayesian network?

  – This time taking into account that probability tables have to sum to 1

| 6 | 12 | 20 | $2^6$-1 |

# Example for BN construction: Fire Diagnosis



- **All variables are Boolean**

- **How many probabilities do we need to specify for this network?**
  - This time taking into account that probability tables have to sum to 1
    - P(Tampering): 1 probability
    - P(Alarm|Tampering, Fire): 4
      1 probability for each of the 4 instantiations of the parents
    - In total: 1+1+4+2+2+2 = 12

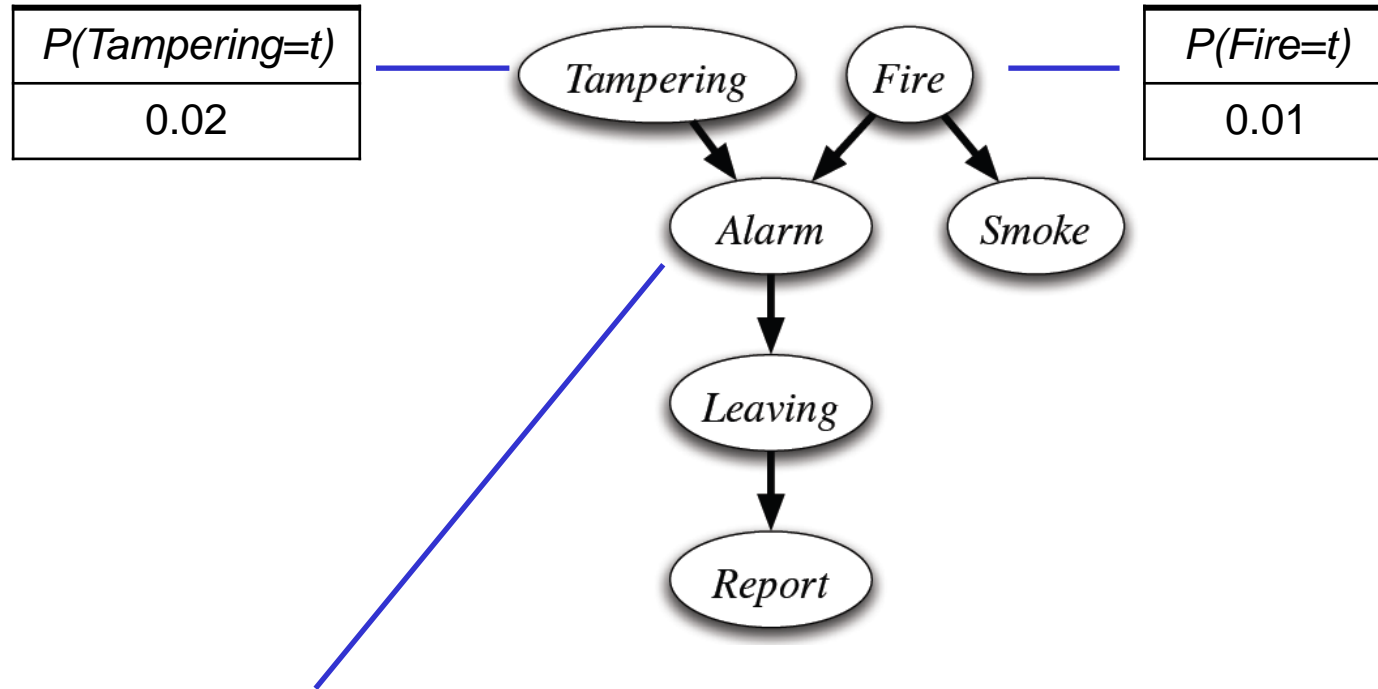# Example for BN construction: Fire Diagnosis



| P(Tampering=t) | P(Tampering=f) |
|:---:|:---:|
| 0.02 | 0.98 |

We don't need to
store P(Tampering=f)
since probabilities sum to 1
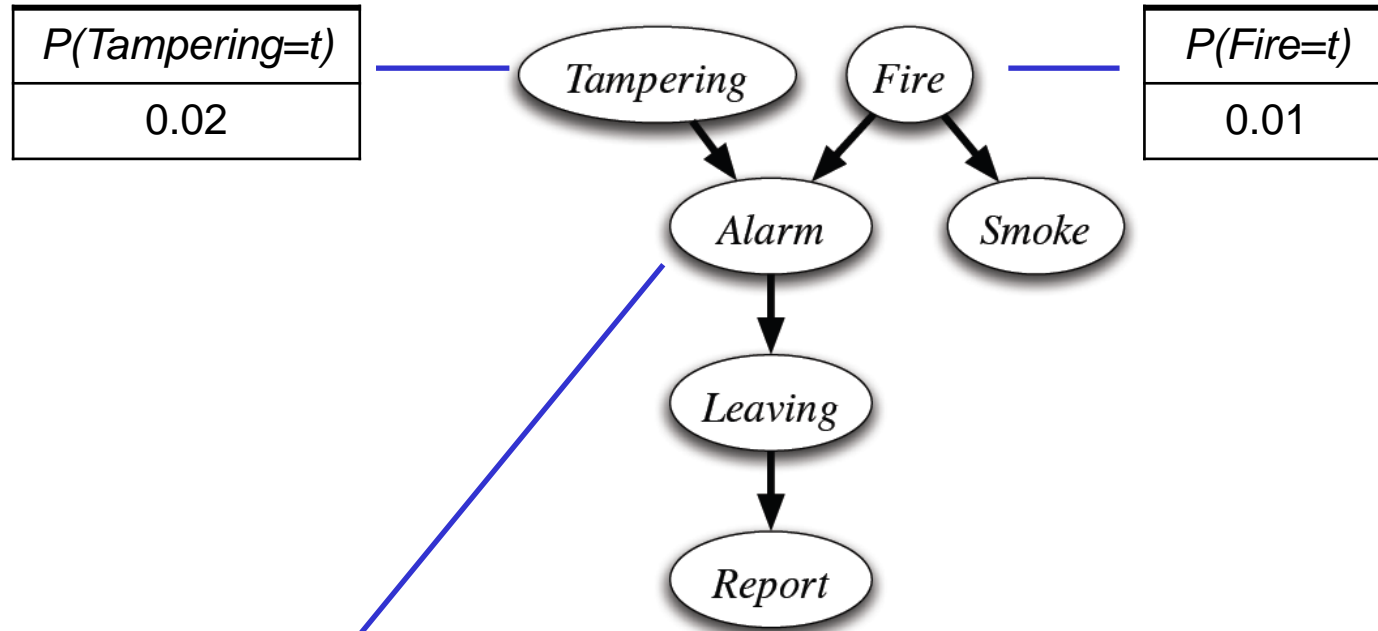
# Example for BN construction: Fire Diagnosis

| P(Tampering=t) |
|:---:|
| 0.02 |

| P(Fire=t) |
|:---:|
| 0.01 |



| Tampering T | Fire F | P(Alarm=t\|T,F) | P(Alarm=f\|T,F) |
|:---:|:---:|:---:|:---:|
| t | t | 0.5 | 0.5 |
| t | f | 0.85 | 0.15 |
| f | t | 0.99 | 0.01 |
| f | f | 0.0001 | 0.9999 |

We don't need to store P(Alarm=f|T,F) since probabilities sum to 1

Each row of this table is a conditional probability distribution

Each column of this table is a conditional probability distribution
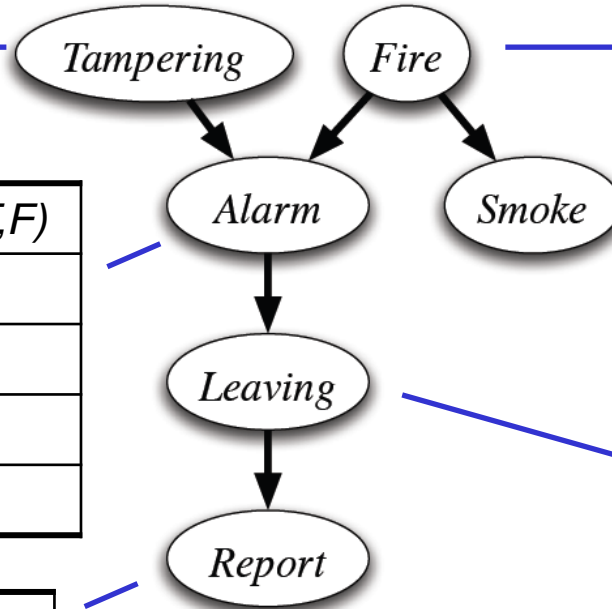
# Example for BN construction: Fire Diagnosis

| P(Tampering=t) |
|:---:|
| 0.02 |

| P(Fire=t) |
|:---:|
| 0.01 |



| Tampering T | Fire F | P(Alarm=t\|T,F) |
|:---:|:---:|:---:|
| t | t | 0.5 |
| t | f | 0.85 |
| f | t | 0.99 |
| f | f | 0.0001 |

We don't need to store P(Alarm=f|T,F) since probabilities sum to 1

Each row of this table is a conditional probability distribution

# Example for BN construction: Fire Diagnosis

| P(Tampering=t) |
|:---:|
| 0.02 |

| P(Fire=t) |
|:---:|
| 0.01 |

*Tampering* → *Alarm*

*Fire* → *Alarm*, *Fire* → *Smoke*

| Tampering T | Fire F | P(Alarm=t|T,F) |
|:---:|:---:|:---:|
| t | t | 0.5 |
| t | f | 0.85 |
| f | t | 0.99 |
| f | f | 0.0001 |

| Fire F | P(Smoke=t |F) |
|:---:|:---:|
| t | 0.9 |
| f | 0.01 |

*Alarm* → *Leaving* → *Report*

| Alarm | P(Leaving=t|A) |
|:---:|:---:|
| t | 0.88 |
| f | 0.001 |

| Leaving | P(Report=t|A) |
|:---:|:---:|
| t | 0.75 |
| f | 0.01 |

P(Tampering=t, Fire=f, Alarm=t, Smoke=f, Leaving=t, Report=t)

# Example for BN construction: Fire Diagnosis

| P(Tampering=t) |
|:---:|
| 0.02 |

| P(Fire=t) |
|:---:|
| 0.01 |

| Tampering T | Fire F | P(Alarm=t|T,F) |
|:---:|:---:|:---:|
| t | t | 0.5 |
| t | f | 0.85 |
| f | t | 0.99 |
| f | f | 0.0001 |

| Fire F | P(Smoke=t |F) |
|:---:|:---:|
| t | 0.9 |
| f | 0.01 |

| Alarm | P(Leaving=t|A) |
|:---:|:---:|
| t | 0.88 |
| f | 0.001 |

| Leaving | P(Report=t|A) |
|:---:|:---:|
| t | 0.75 |
| f | 0.01 |

P(Tampering=t, Fire=f, Alarm=t, Smoke=f, Leaving=t, Report=t)

= P(Tampering=t) × P(Fire=f) × P(Alarm=t|Tampering=t,Fire=f)
    × P(Smoke=f|Fire=f) × P(Leaving=t|Alarm=t)
    × P(Report=t|Leaving=t)

= 0.02 × (1-0.01) × 0.85 × (1-0.01) × 0.88 × 0.75

# What if we use a different ordering?

- Important for assignment 4, question 2:
- Say, we use the following order:
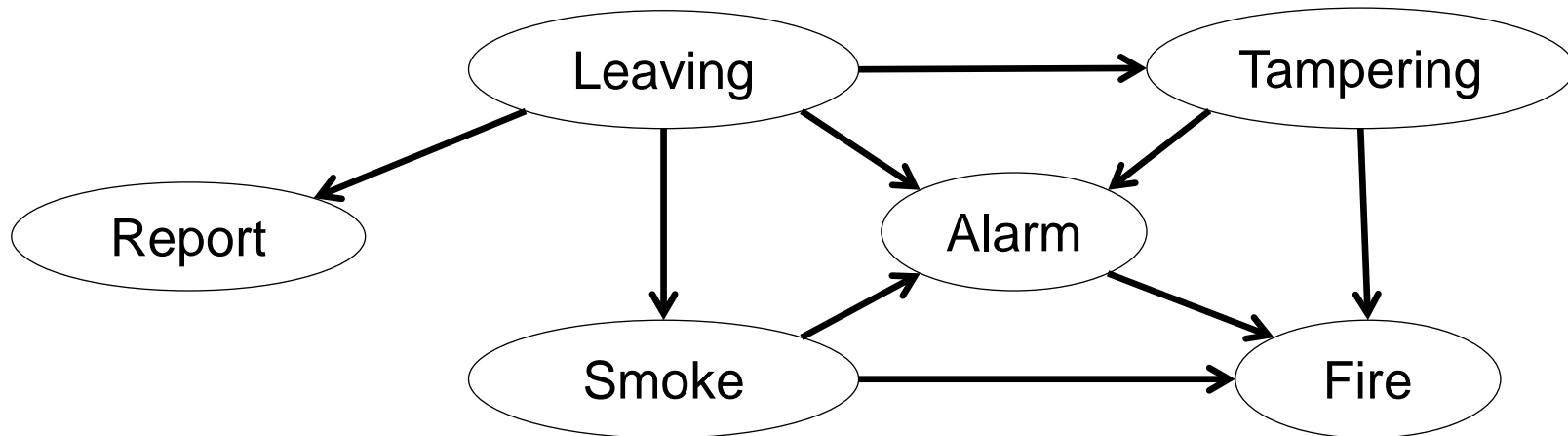  - Leaving; Tampering; Report; Smoke; Alarm; Fire.



- We end up with a completely different network structure!
- Which of the two structures is better (think computationally)?

The previous structure

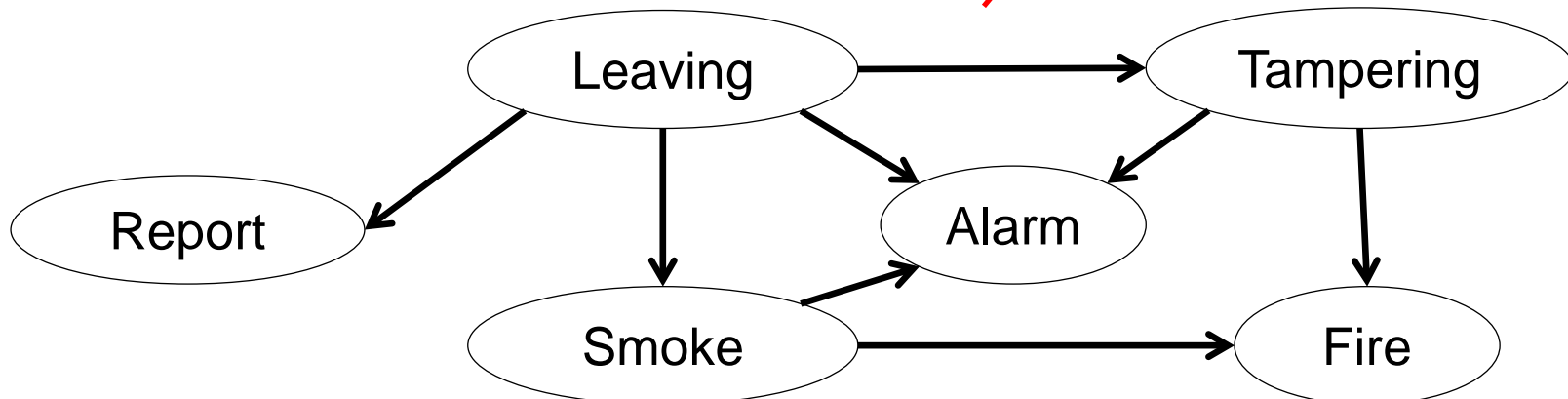This structure

# What if we use a different ordering?

- Important for assignment 4, question 2:
- Say, we use the following order:
  - Leaving; Tampering; Report; Smoke; Alarm; Fire.



- We end up with a completely different network structure!
- Which of the two structures is better (think computationally)?
  - In the last network, we had to specify 12 probabilities
  - Here? 1 + 2 + 2 + 2 + 8 + 8 = 23
  - The causal structure typically leads to the most compact network
    - Compactness typically enables more efficient reasoning

25

# Are there wrong network structures?

- Important for assignment 4, question 2
- Some variable orderings yield more compact, some less compact structures
    - Compact ones are better
    - But all representations resulting from this process are correct
    - One extreme: the fully connected network is always correct but rarely the best choice
- How can a network structure be wrong?
    - If it misses directed edges that are required
    - E.g. an edge is missing below: Fire ⊥̸ Alarm | {Tampering, Smoke}
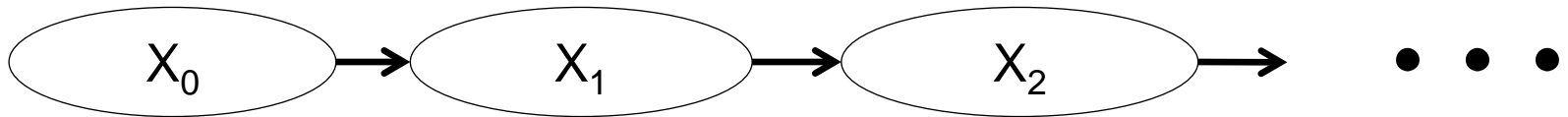
# Lecture Overview

- Recap: marginal and conditional independence

- Bayesian Networks Introduction
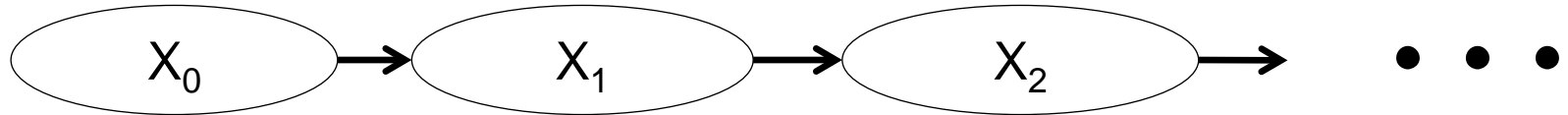
Hidden Markov Models
- Rainbow Robot Example

# Markov Chains

- A Markov chain is a special kind of belief network:

$$X_0 \longrightarrow X_1 \longrightarrow X_2 \longrightarrow \quad \bullet\ \bullet\ \bullet$$

- $X_t$ represents a state at time t.

- Its dependence structure yields: $P(X_t|X_1, \ldots, X_{t-1}) = P(X_t|X_{t-1})$
  - This conditional probability distribution is called the state transition probability
  - Intuitively $X_t$ conveys all of the information about the history that can affect the future states:
    "The past is independent of the future given the present."

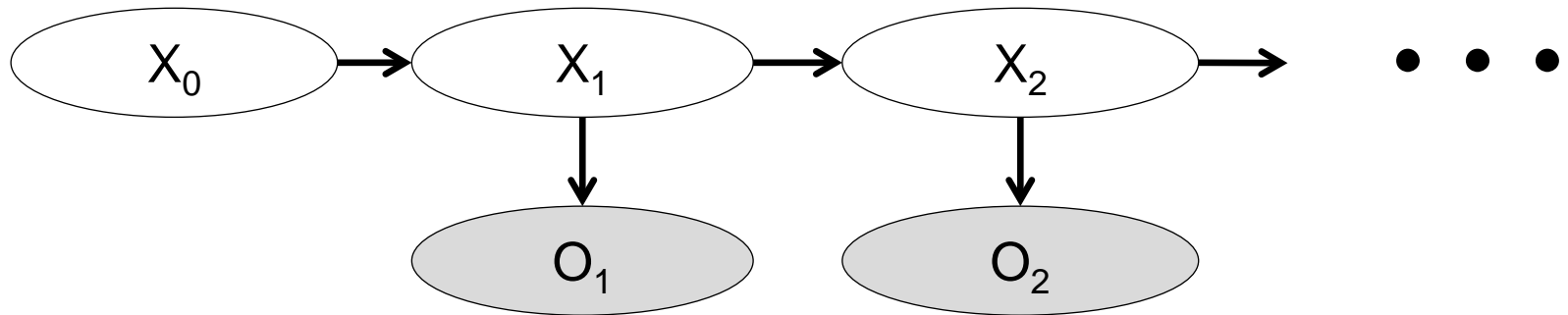- JPD of a Markov Chain: $P(X_0, \ldots, X_T) = P(X_0) \times \prod_{t=1}^{T} P(X_t|X_{t-1})$

# Stationary Markov Chains

$$X_0 \longrightarrow X_1 \longrightarrow X_2 \longrightarrow \quad \bullet \; \bullet \; \bullet$$

- A stationary Markov chain is when
  - All state transition probability tables are the same
  - I.e., for all $t > 0$, $t' > 0$: $P(X_t|X_{t-1}) = P(X_{t'}|X_{t'-1})$

- We only need to specify $P(X_0)$ and $P(X_t|X_{t-1})$.
  - Simple model, easy to specify
  - Often the natural model
  - The network can extend indefinitely
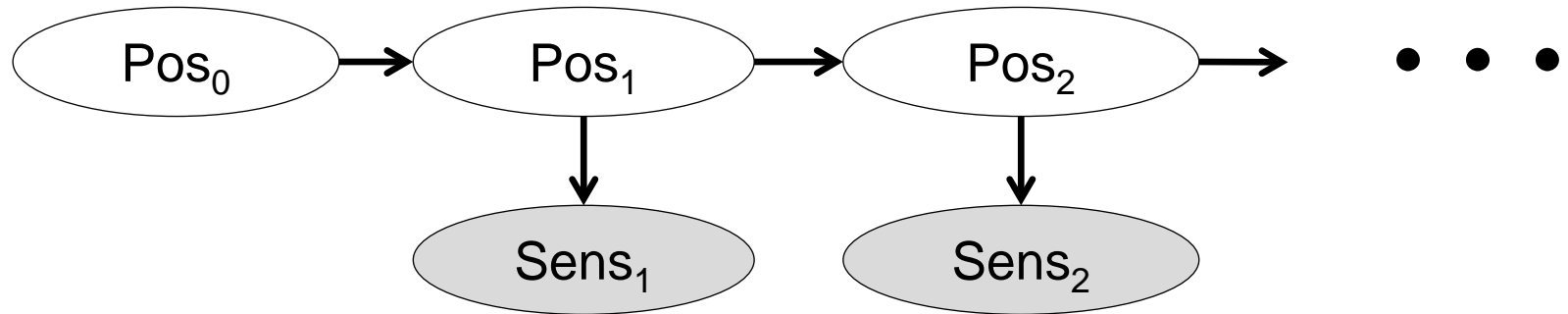
# Hidden Markov Models (HMMs)

- A Hidden Markov Model (HMM) is a Markov chain plus a noisy observation about the state at each time step:

```
( X_0 ) → ( X_1 ) → ( X_2 ) →   • • •
              ↓         ↓
           ( O_1 )   ( O_2 )
```

- Same conditional probability tables at each time step
  - The state transition probability $P(X_t|X_{t-1})$
    - also called the system dynamics
  - The observation probability $P(O_t|X_t)$
    - also called the sensor model

- JPD of an HMM: $P(X_0, \ldots, X_T, O_1, \ldots, O_T)$
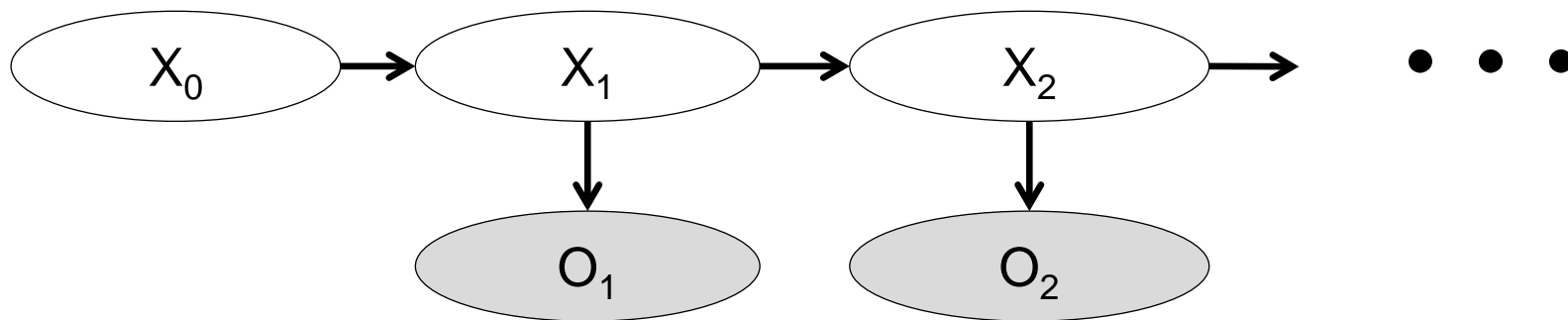  $= P(X_0) \times \prod_{t=1}^{T} P(X_t|X_{t-1}) \times \prod_{t=1}^{T} P(O_t|X_{t-1})$

# Example HMM: Robot Tracking

- Robot tracking as an HMM:



- Robot is moving at random: $P(Pos_t|Pos_{t-1})$
- Sensor observations of the current state $P(Sens_t|Pos_t)$

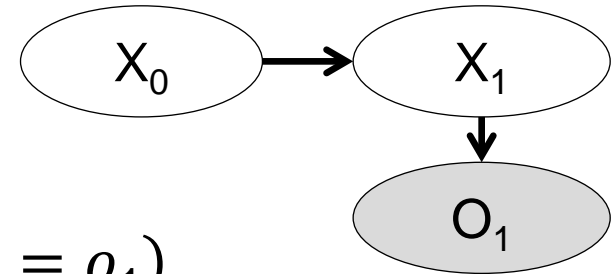# Filtering in Hidden Markov Models (HMMs)



- **Filtering** problem in HMMs:
  at time step t, we would like to know $P(X_t | O_1, \ldots, O_t)$

- We will derive simple update equations:
  - Compute $P(X_t | O_1, \ldots, O_t)$ if we already know $P(X_{t-1} | O_1, \ldots, O_{t-1})$

# HMM Filtering: first time step

By applying marginalization over $X_0$ "backwards":

$P(X_1 | O_1 = o_1)$



$= \sum_{x \in dom(X_0)} P(X_1, X_0 = x | O_1 = o_1)$

Direct application of Bayes rule

$= \sum_{x \in dom(X_0)} \frac{P(O_1 = o_1 | X_1, X_0 = x) \times P(X_1, X_0 = x)}{P(O_1 = o_1)}$

$O_1 \perp\!\!\!\perp X_0 | X_1$ and product rule

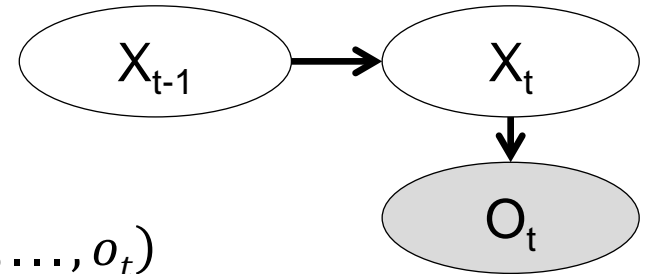$= \sum_{x \in dom(X_0)} \frac{P(O_1 = o_1 | X_1) \times P(X_1 | X_0 = x) \times P(X_0 = x)}{P(O_1 = o_1)}$

Normalize to make the probability to sum to 1.

$\propto \sum_{x \in dom(X_0)} P(O_1 = o_1 | X_1) \times P(X_1 | X_0 = x) \times P(X_0 = x)$

# HMM Filtering: general time step t

$$P(X_t | o_1, \ldots, o_t)$$

$$= \sum_{x \in dom(X_{t-1})} P(X_t, X_{t-1} = x | o_1, \ldots, o_t)$$

Direct application of Bayes rule

$$= \sum_{x \in dom(X_{t-1})} \frac{P(o_t | X_t, X_{t-1} = x, o_1, \ldots, o_{t-1}) \times P(X_t, X_{t-1} = x | o_1, \ldots, o_{t-1})}{P(o_t)}$$

$O_t \perp\!\!\!\perp \{X_{t-1}, O_1, \ldots, O_{t-1}\} \mid X_t$ and $X_t \perp\!\!\!\perp \{O_1, \ldots, O_{t-1}\} \mid X_{t-1}$

$$= \sum_{x \in dom(X_{t-1})} \frac{P(o_t | X_t) \times P(X_t | X_{t-1} = x) \times P(X_{t-1} = x | o_1, \ldots, o_{t-1})}{P(o_t)}$$

Normalize to make the probability to sum to 1.

$$\propto \sum_{x \in dom(X_{t-1})} P(o_t | X_t) \times P(X_t | X_{t-1} = x) \times P(X_{t-1} = x | o_1, \ldots, o_{t-1})$$

$X_{t-1} \longrightarrow X_t$

$O_t$

# HMM Filtering Summary

- Initialize belief state at time 0: $P(X_0)$
  - In Rainbow Robots, we initialize this for you: $P(Pos_t)$

- At each time step, update belief state given new observation:

$$P(X_t | o_1, \ldots, o_t)$$

$$\propto \sum_{x \in dom(X_{t-1})} P(o_t | X_t) \times P(X_t | X_{t-1} = x) \times P(X_{t-1} = x | o_1, \ldots, o_{t-1})$$

Observation probability

Transition probability

We already know this from the previous step

- Rainbow Robot example
  - take the last belief state,
  - multiply it with the transition probability $P(Pos_t | Pos_{t-1})$
  - multiply it with the observation probability $P(Sens_t | Pos_t)$
  - and normalize

35

# Learning Goals For Today's Class

- Build a Bayesian Network for a given domain
- Compute the representational savings in terms of number of probabilities required

---

- Assignment 4 available on WebCT
  - Due Monday, April 4
    - Can only use 2 late days
    - So we can give out solutions to study for the final exam
  - Final exam: Monday, April 11
    - Less than 3 weeks from now

  - You should now be able to solve questions 1, 2, and 5
    - Material for Question 3: Friday, and wrap-up on Monday
    - Material for Question 4: next week