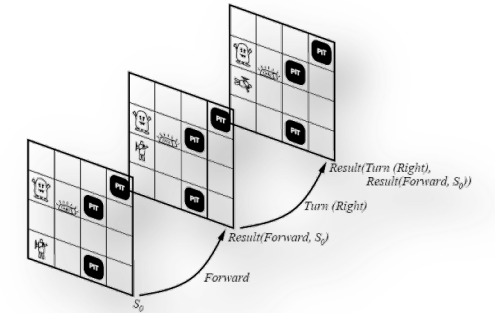# Introduction to Artificial Intelligence

**Roman Barták**

Department of Theoretical Computer Science and Mathematical Logic

In **situation calculus**, we view the world as a series of snapshots (**time slices**). A similar approach can be applied in **probabilistic reasoning about time**.

Each time slice (**state**) is described as a set of random variables:

- **hidden** (not observable) random **variables $X_t$**
  *describe the actual state*
- **observable** random **variables $E_t$** (with observed values $e_t$)
  *describe what we observe about the state*
- **t** is an identification of the time slice (we assume **discrete time** with uniform time steps)

*Notation:*

$X_{a:b}$ denotes a set of variables from $X_a$ to $X_b$

We need to describe evolution of states and how observations depend on states.

## Transition model

specifies the probability distribution over the latest state variables given the previous values $P(X_t \mid X_{0:t-1})$

*Simplifying assumptions:*

— state depends on previous state only (**Markov assumption**): $P(X_t \mid X_{0:t-1}) = P(X_t \mid X_{t-1})$

— all transitions tables $P(X_t \mid X_{t-1})$ are identical for all t (**stationary process**)

## Sensor (observation) model

describes how the evidence (observed) variables $E_t$ depend on other variables $P(E_t \mid X_{0:t}, E_{1:t-1})$

*Simplifying assumption:*

— observation depends on current state only (**sensor Markov assumption**): $P(E_t \mid X_{0:t}, E_{1:t-1}) = P(E_t \mid X_t)$

You are the security guard stationed at a secret underground installation and you want to know whether it is **raining today**:
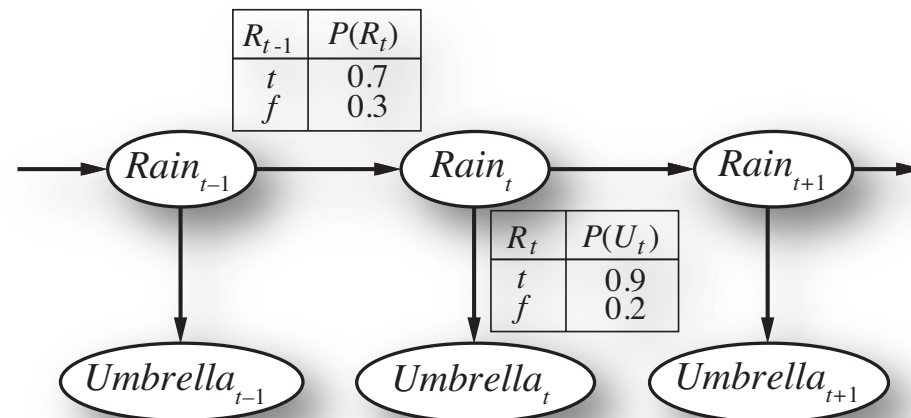
- hidden random variable $R_t$

But your only access to the outside world occurs each morning when you see the the director coming in **with, or without, an umbrella**.

- observable random variable $U_t$

The transition and sensor models can be described using a **Bayesian network**.

In addition to $P(X_t \mid X_{t-1})$ and $P(E_t \mid X_t)$ we need to say how everything gets started $P(X_0)$ (= $\langle 0.5, 0.5 \rangle$, for example).

| $R_{t-1}$ | $P(R_t)$ |
|---|---|
| $t$ | 0.7 |
| $f$ | 0.3 |

$Rain_{t-1}$ → $Rain_t$ → $Rain_{t+1}$

| $R_t$ | $P(U_t)$ |
|---|---|
| $t$ | 0.9 |
| $f$ | 0.2 |

$Umbrella_{t-1}$  $Umbrella_t$  $Umbrella_{t+1}$

We have a specification of the complete joint distribution:

$$P(X_{0:t}, E_{1:t}) = P(X_0) \prod_i P(X_i \mid X_{i-1}) \, P(E_i \mid X_i)$$

**Filtering**: the task of computing the posterior distribution over the *most recent state,* given all evidence to date
$P(X_t \mid e_{1:t})$

Where am I now?

**Prediction**: the task of computing the posterior distribution over the *future state*, given all evidence to date
$P(X_{t+k} \mid e_{1:t})$ for k>0

Where will I be in future?

**Smoothing**: the task of computing posterior distribution over a *past state*, given all evidence up to the present
$P(X_k \mid e_{1:t})$ for k: $0 \leq k < t$

Where was I in past?

**Most likely explanation**: the task to find the sequence of states that is most likely generated a given sequence of observations

What path did I go through?

$\text{argmax}_{x_{1:t}} P(x_{1:t} \mid e_{1:t})$

The task of computing the posterior distribution over the *most recent state,* given all evidence to date – $P(X_t|e_{1:t})$.

A useful filtering algorithm needs to maintain a current state estimate and update it, rather than going back over (**recursive estimation**):

$P(X_{t+1}|e_{1:t+1}) = f(e_{t+1}, P(X_t|e_{1:t}))$

How to define the function f?

$P(X_{t+1}|e_{1:t+1}) = P(X_{t+1}|e_{1:t}, e_{t+1})$

$= \alpha \, P(e_{t+1}|X_{t+1}, e_{1:t}) \, P(X_{t+1}|e_{1:t})$

$= \alpha \, P(e_{t+1}|X_{t+1}) \, P(X_{t+1}|e_{1:t})$

$= \alpha \, P(e_{t+1}|X_{t+1}) \sum_{x_t} P(X_{t+1}|x_t, e_{1:t}) \, P(x_t|e_{1:t})$

$= \alpha \, P(e_{t+1}|X_{t+1}) \sum_{x_t} P(X_{t+1}|x_t) \, P(x_t|e_{1:t})$

> **Bayes rule**

> **sensor Markov assumption**

> **conditioning**

A message $f_{1:t}$ is propagated forward over the sequence:

$P(X_t|e_{1:t}) = f_{1:t}$

$f_{1:t+1} = \alpha \, \text{FORWARD}(f_{1:t}, e_{t+1})$

$f_{1:0} = P(X_0)$

The task of computing the posterior distribution over the *future state*, given all evidence to date – $P(\mathbf{X}_{t+k} \mid \mathbf{e}_{1:t})$ for some k>0.

We can see this task as filtering without the addition of new evidence:

$$P(\mathbf{X}_{t+k+1} \mid \mathbf{e}_{1:t}) = \Sigma_{\mathbf{x}_{t+k}} \; P(\mathbf{X}_{t+k+1} \mid \mathbf{x}_{t+k}) \; P(\mathbf{x}_{t+k} \mid \mathbf{e}_{1:t})$$

After some time (**mixing time**) the predicted distribution converges to the **stationary distribution** of the Markov process and remains constant.

The task of computing posterior distribution over a *past state*, given all evidence up to the present – $P(X_k|e_{1:t})$ for k: $0 \le k < t$.
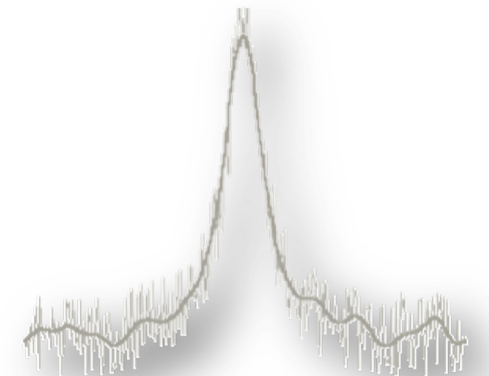
We again exploit a recursive message-passing approach, now in two directions.

$P(X_k|e_{1:t}) = P(X_k|e_{1:k},e_{k+1:t})$      **Bayes rule**

     $= \alpha \; P(X_k|e_{1:k}) \; P(e_{k+1:t}|X_k,e_{1:k})$

     $= \alpha \; P(X_k|e_{1:k}) \; P(e_{k+1:t}|X_k)$      **conditional independence**

     $= \alpha \; f_{1:k} \times b_{k+1:t}$

$P(e_{k+1:t}|X_k) = \sum_{x_{k+1}} P(e_{k+1:t}|X_k,x_{k+1}) \; P(x_{k+1}|X_k)$      **conditioning**

     $= \sum_{x_{k+1}} P(e_{k+1:t}|x_{k+1}) \; P(x_{k+1}|X_k)$      **conditional independence**

     $= \sum_{x_{k+1}} P(e_{k+1},e_{k+2:t}|x_{k+1}) \; P(x_{k+1}|X_k)$

     $= \sum_{x_{k+1}} P(e_{k+1}|x_{k+1}) \; P(e_{k+2:t}|x_{k+1}) \; P(x_{k+1}|X_k)$      **conditional independence**

Using the backward message-passing notation:

     $P(e_{k+1:t}|X_k) = b_{k+1:t}$

     $b_{k+1:t} = \text{BACKWARD}(b_{k+2:t}, e_{k+1})$

     $b_{t+1:t} = P(e_{t+1:t}|X_t) = P(|X_t) = 1$

The task to find the *sequence of states that is most likely* generated a given sequence of observations – $\text{argmax}_{x_{1:t}} P(x_{1:t} \mid e_{1:t})$.

Note: This is different from smoothing for each past state and taking the sequence of most probable states!

We can see each sequence as a **path through a graph** whose nodes are possible states at each time step.
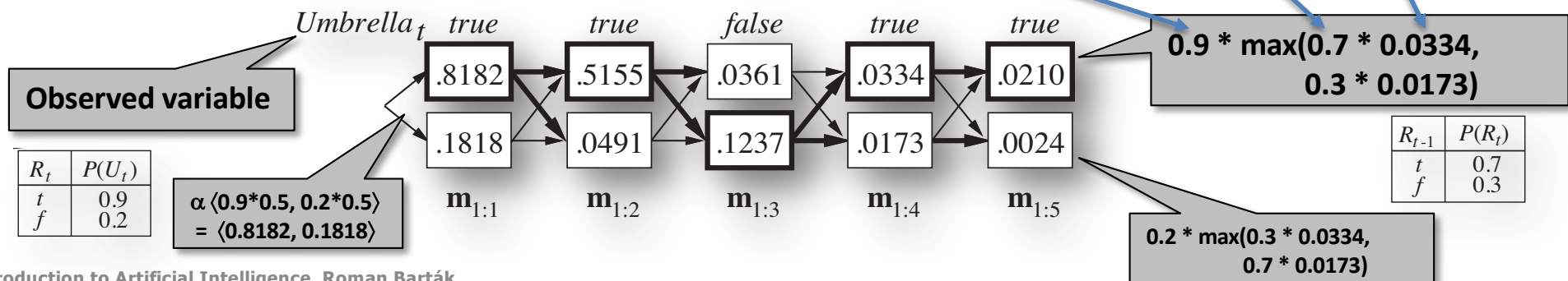


The most likely path to a given state consists of the most likely path to some previous state followed by a transition to that state.

This can be described using a **recursive formula** (**Viterbi algorithm**):

$$\max_{x_1,\ldots,x_t} P(x_1,\ldots,x_t,X_{t+1} \mid e_{1:t+1})$$
$$= \alpha\, P(e_{t+1} \mid X_{t+1}) \max_{x_t} (P(X_{t+1} \mid x_t) \max_{x_1,\ldots,x_{t-1}} P(x_1,\ldots,x_t \mid e_{1:t}))$$

Again, we use an approach of forward message passing:

$$m_{1:t} = \max_{x_1,\ldots,x_{t-1}} P(x_1,\ldots,x_{t-1},X_t \mid e_{1:t}), \qquad m_{1:t+1} = P(e_{t+1} \mid X_{t+1}) \max_{x_t} (P(X_{t+1} \mid x_t)\, m_{1:t})$$



**Observed variable**

| $R_t$ | $P(U_t)$ |
|---|---|
| t | 0.9 |
| f | 0.2 |

$\alpha \langle 0.9*0.5,\ 0.2*0.5 \rangle$
$= \langle 0.8182,\ 0.1818 \rangle$

$Umbrella_t$ true, true, false, true, true

.8182 → .5155 → .0361 → .0334 → .0210
.1818 → .0491 → .1237 → .0173 → .0024

$m_{1:1}$  $m_{1:2}$  $m_{1:3}$  $m_{1:4}$  $m_{1:5}$

0.9 * max(0.7 * 0.0334, 0.3 * 0.0173)

0.2 * max(0.3 * 0.0334, 0.7 * 0.0173)

| $R_{t-1}$ | $P(R_t)$ |
|---|---|
| t | 0.7 |
| f | 0.3 |

Assume that the state of process is described by a single discrete random variable $X_t$ (there is also a single evidence variable $E_t$).

This is called a **hidden Markov model** (HMM).

This restricted model allows for a simple and elegant **matrix implementation** of all the basic algorithms.

Assume that variable $X_t$ takes values from the set $\{1,...S\}$, where S is the number of possible states.

The **transition model** $\mathbf{P}(X_t \mid X_{t-1})$ becomes an S×S matrix $\mathbf{T}$, where:

$$\mathbf{T}_{(i,j)} = P(X_t = j \mid X_{t-1}=i)$$

We also put the **sensor model** in matrix form. Now we know the value of the evidence variable $e_t$ so we describe $P(E_t = e_t \mid X_t=i)$, using a diagonal matrix $\mathbf{O}_t$, where:

$$\mathbf{O}_{t\,(i,i)} = P(E_t = e_t \mid X_t=i)$$

**The forward message propagation** (from filtering)

$P(X_t | e_{1:t}) = f_{1:t}$

$f_{1:t+1} = \alpha\, P(e_{t+1} | X_{t+1})\, \Sigma_{x_t}\, P(X_{t+1} | x_t)\, P(x_t | e_{1:t})$

can be reformulated using matrix operations (message $f_{1:t}$ is modelled as a one-column matrix) as follows:

$T_{(i,j)} = P(X_t = j\ |\ X_{t-1} = i)$

$O_{t\,(i,i)} = P(E_t = e_t\ |\ X_t = i)$

$\textcolor{red}{f_{1:t+1} = \alpha\, O_{t+1}\, T^\top\, f_{1:t}}$

**The backward message propagation** (from smoothing)

$P(e_{k+1:t} | X_k) = b_{k+1:t}$

$b_{k+1:t} = \Sigma_{x_{k+1}}\, P(e_{k+1} | x_{k+1})\, P(e_{k+2:t} | x_{k+1})\, P(x_{k+1} | X_k)$

can be reformulated using matrix operations (message $b_{k:t}$ is modelled as a one-column matrix) as follows:

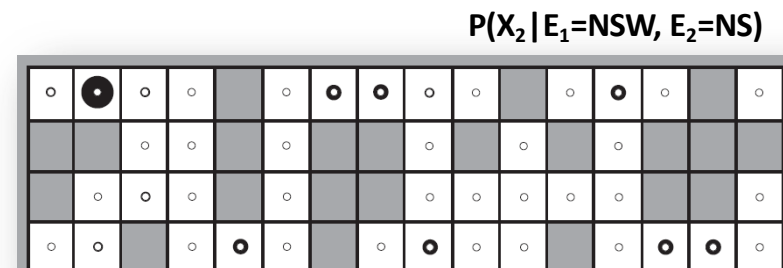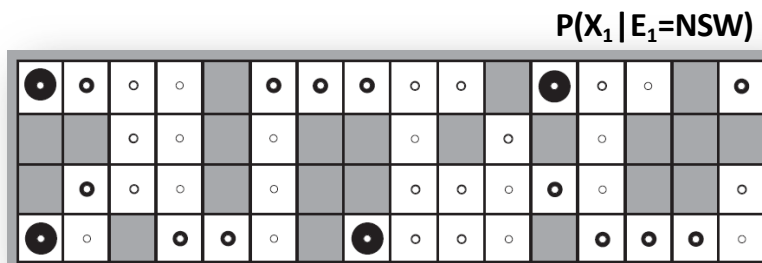$\textcolor{red}{b_{k+1:t} = T\, O_{k+1}\, b_{k+2:t}}$

Assume a robot that moves randomly in a grid world, has a map of the world and (noisy) sensors reporting obstacles laying immediately to the north, south, east, and west. The robot needs to find its location.

## A possible model:

- **random variables** $X_t$ describe robot's location at times t
  - possible values are 1,..,n for n locations
  - Nb(i) – a set of neighboring locations for location i
- **transition tables** (random move)
  - $P(X_{t+1}=j|X_t=i) = 1/|Nb(i)|$,     if $j \in Nb(i)$,
                    0,                otherwise
- **sensor variables** $E_t$ describe observations (evidence) at times t (four sensor for four directions NSEW)
  - values indicate detection of obstacle at a given direction NSEW (16 values for all directions)
  - assume that sensor's error rate is $\varepsilon$
- **sensor tables**
  - $P(E_t=e_t|X_t=i) = (1-\varepsilon)^{4-d_{it}} \varepsilon^{d_{it}}$
    where $d_{it}$ is the number of deviations of observation $e_t$ from the true values for square i

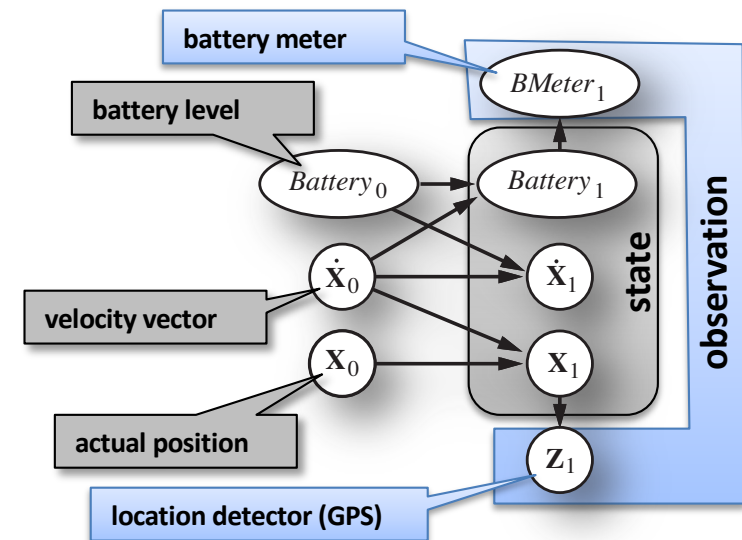$P(X_1|E_1=NSW)$



$P(X_2|E_1=NSW, E_2=NS)$

# **Dynamic Bayesian network** (DBN) is a Bayesian network that represents a temporal probability model.

the variables and links are exactly replicated from slice to slice
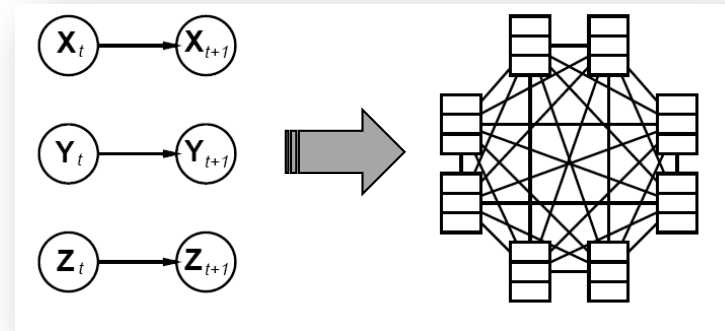
It is enough to describe one slice.

- prior distribution $P(X_0)$
- transition model $P(X_1 | X_0)$
- sensor model $P(E_1 | X_1)$



Each state variable has parents either at the same slice or in the previous slice (Markov assumption).

A hidden Markov model is a special case of a dynamic Bayesian network.

Similarly, a dynamic Bayesian network can be encoded as a hidden Markov model



> one random variable in HMM whose
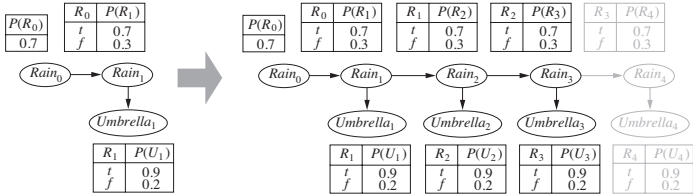> values are n-tuples of values
> of state variables in DBN

**What is the difference?**

> The relationship between DBN and HMM is roughly analogous to the relationship between ordinary Bayesian networks and full tabulated joint distribution.
>
> — DBN with 20 Boolean state variables, each of which has three parents
>   - the transition model has $20 \times 2^3 = 160$ probabilities
> — Hidden Markov model has one random variable with $2^{20}$ values
>   - the transition model has $2^{20} \times 2^{20} \approx 10^{12}$ probabilities
>   - HMM requires much more space and inference is much more expensive

Dynamic Bayesian networks are Bayesian networks and we already have algorithms for inference in Bayesian networks.

We can construct the full Bayesian network representation of a DBN by replicating slices to accommodate the observations (**unrolling**).
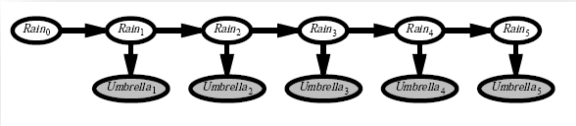
## Exact inference:

If applied naively, its complexity will increase with time (due to more slices).

We can use **variable elimination** and keep in memory only last two slices (via summing out the variables from previous slices).

The bad news are that "constant" space to represent the **largest factor will be exponential** in the number of state variables.

## Approximate inference:

We sample non-evidence nodes of the network in topological order, weighting each sample by the likelihood in accords to the observed evidence variables.

But samples are generated completely independently of the evidence!

Hence, the weights of samples will decrease so to keep accuracy we need to **increase the number of samples exponentially** with t.

We can exploit probability theory when reasoning about time. Specifically, when **transitions are uncertain,** and environment is **partially observable** via sensors.

We use **transition and observation models** with **Markov assumptions**.

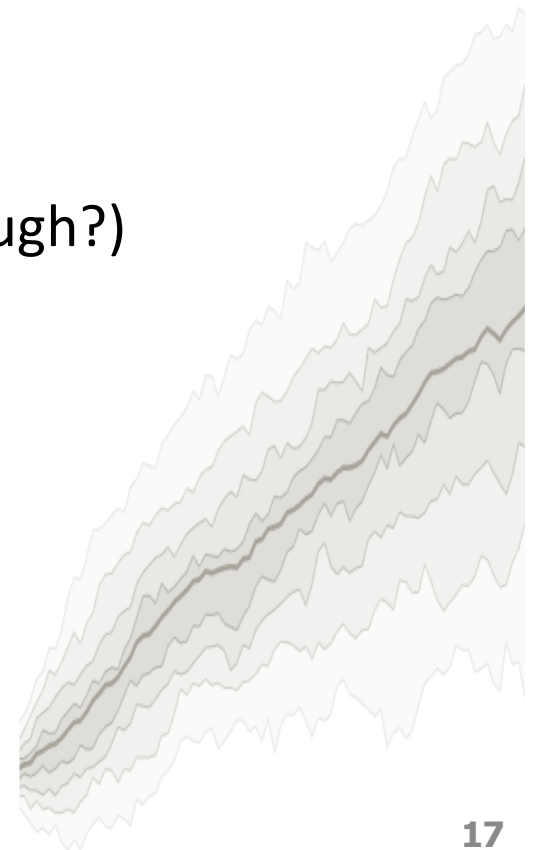**Basic inference tasks** (exploit recursive formulas):
- **filtering** (where am I now?)
- **prediction** (where will I be in future?)
- **smoothing** (where was I in past?)
- **most likely explanation** (what path did I go through?)

**Hidden Markov Model**
- one state variable and one observation variable
- simplified inference using matrix operations

**Dynamic Bayesian Network**
- compact representation via more CPTs