# AI Fundamentals: uncertain reasoning

*Maria Simi*

# Probabilistic Reasoning over time
# AIMA chapter 15

LESSON 3: TIME AND UNCERTAINTY - INFERENCE IN TEMPORAL MODELS

# Reasoning in a changing world

So far, we have been doing uncertain reasoning in a static world.

For reasoning in an evolving world an agent needs:

- A **belief state:** the states of the world that are possible
- A **transition model:** to predict how the world will evolve
- A **sensor model:** to update the belief state from perceptions

A changing world is modeled using a variable for each aspect of the world state *at each point in time* (fluents). We use probabilities to quantify how likely a world is.

The transition and sensor model themselves may be uncertain:

- the *transition model* gives the probability distribution of the variables at time $t$, given the state of the world at past times;
- the *sensor model* describes the probability of each percept at time $t$, given the current state of the world.

# Examples

1. **Treating a diabetic patient.** The task is to assess the current state of the patient, including the actual blood sugar level and insulin level from other observable parameters such as recent insulin doses, food intake, blood sugar measurements, and other physical signs.

   Dynamic aspects are essential for a correct prediction and prescription of insulin and food, in this case as in many other real-world scenarios.

2. **The umbrella example, our toy example.**

   *You are the security guard stationed at a secret underground installation. You want to know whether it's raining today, but your only access to the outside world occurs each morning when you see the director coming in with, or without, an umbrella.*

# States and observations

We view the world as a series of *snapshots*, or **time slices**, each of which contains a set of random variables, some observable and some not.

$X_t$: will denote the set of **state variables**, unobservable (hidden) at time $t$

$E_t = e_t$ will denote the set of observations (**evidence variables**) at time $t$; $e_t$ their values

We will assume that the state sequence starts at $t=0$. The distance between time slices is fixed.

We will use the notation $a{:}b$ to denote the sequence of integers from $a$ to $b$ (inclusive), and the notation $X_{a{:}b}$ to denote the set of variables from $X_a$ to $X_b$

Our umbrella world is represented by a sequence of evidence variables $E_t = \{U_t\}$ (whether the umbrella appears at time $t$) and state unobservable variables $X_t = \{R_t\}$ (it's raining):

state variables: $R_0, R_1, R_2, \ldots$

evidence variables: $U_1, U_2, \ldots$

**Note**: state variable start at 0, evidence variables at time 1.

# Transition model and simplifying assumptions

The transition model specifies **how the world evolves**, i.e. the probability distribution over the latest state variables, given the previous values starting from time 0:

$\mathbf{P}(\mathbf{X}_t \mid \mathbf{X}_{0:t-1})$ $\qquad\qquad\qquad$ $\mathbf{X}_{0:t-1} = X_0, X_1, \ldots X_{t-1}$

The sequence of states can become very large, unbounded as $t$ increases.

**Markov assumptions**: the transition model specifies the probability distribution over the latest state variables, given a **finite fixed number** of previous states:

$\mathbf{P}(\mathbf{X}_t \mid \mathbf{X}_{0:t-1}) = \mathbf{P}(\mathbf{X}_t \mid \mathbf{X}_{t-1})$ $\qquad\qquad$ *first order* Markov process/chain

$\mathbf{P}(\mathbf{X}_t \mid \mathbf{X}_{0:t-1}) = \mathbf{P}(\mathbf{X}_t \mid \mathbf{X}_{t-1}, \mathbf{X}_{t-2})$ $\qquad\quad$ *second order* Markov process/chain
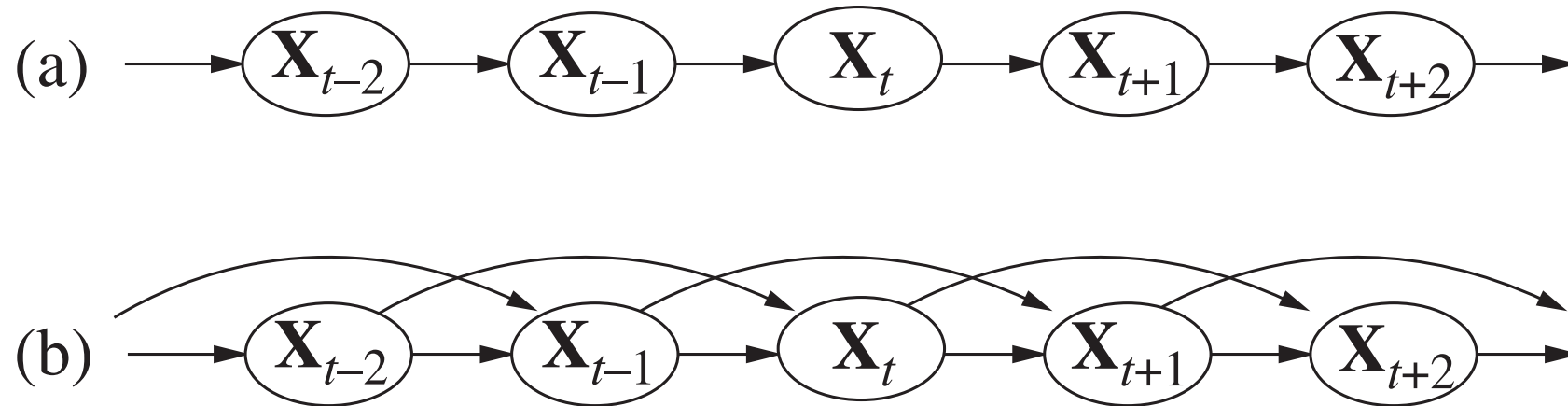
Additionally, we assume a **stationary process**, i.e. the conditional probability distribution is the same for all $t$.

$\mathbf{P}(\mathbf{X}_t \mid \mathbf{X}_{t-1})$ is the same for all $t$ $\qquad\qquad$ $\mathbf{P}(R_t \mid r_{t-1}) = \langle 0.7, 0.3 \rangle$

Change is governed by laws that do not themselves change over time.

# Transition model - graphically

Bayes networks under different Markov [independence] assumptions

(a) $\longrightarrow$ $X_{t-2}$ $\longrightarrow$ $X_{t-1}$ $\longrightarrow$ $X_t$ $\longrightarrow$ $X_{t+1}$ $\longrightarrow$ $X_{t+2}$ $\longrightarrow$

(b) $\longrightarrow$ $X_{t-2}$ $\longrightarrow$ $X_{t-1}$ $\longrightarrow$ $X_t$ $\longrightarrow$ $X_{t+1}$ $\longrightarrow$ $X_{t+2}$ $\longrightarrow$

a.   *first order* Markov process/chain

b.   *second order* Markov process/chain

AI FUNDAMENTALS - M. SIMI

# Sensor model assumption

The sensor/observation model, under **Markov sensor assumption**, postulates that evidence only depends on the current state

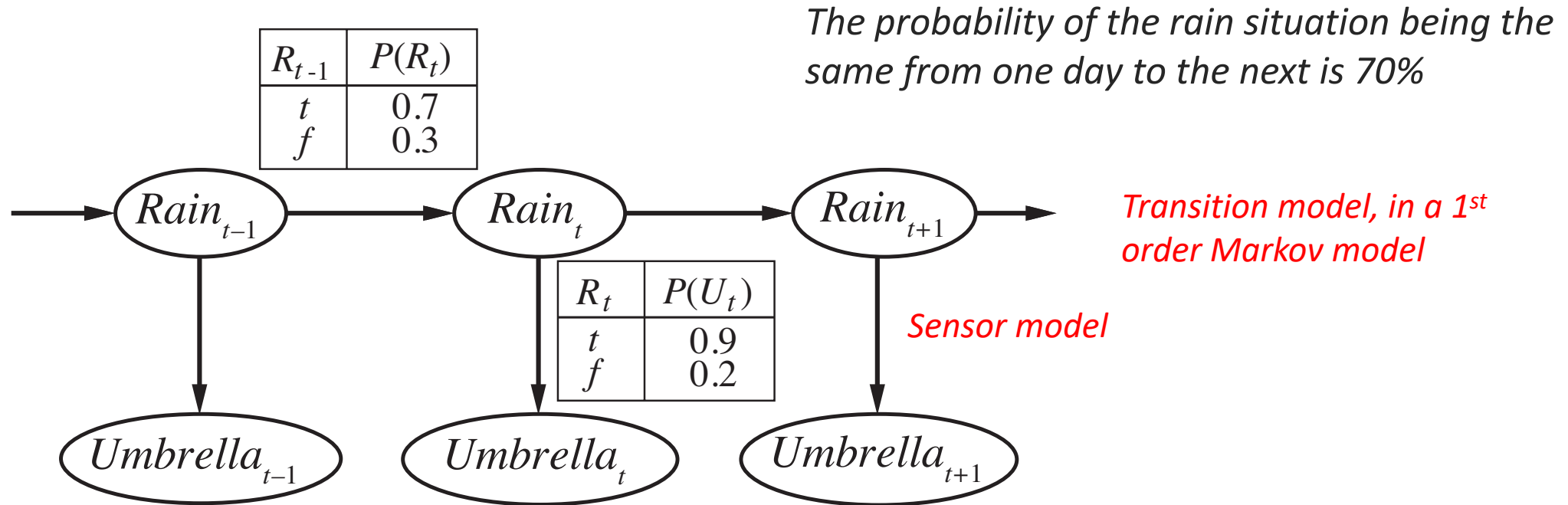$$P(\mathbf{E}_t \mid \mathbf{X}_{0:t}, \mathbf{E}_{0:t-1}) = P(\mathbf{E}_t \mid \mathbf{X}_t)$$

This is a reasonable assumption to make, given the availability of sensors.

**Releasing the assumptions**:

Assuming that 'Rain' only depends on rain the previous day may be a too strong assumption. There are two ways to improve the accuracy of the approximation:

1. Increasing the order of the Markov process model, for example using a second order assumption.
2. Increasing the set of state variables: we could add *Season, Temperature, Humidity, Pressure* … as state variables. This may imply more computation for predicting state variables or adding new sensors.

# The Bayesian net for the *Umbrella* example

*The probability of the rain situation being the same from one day to the next is 70%*

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

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

*Transition model, in a 1st order Markov model*

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

*Sensor model*

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

The transition model is $P(R_t | R_{t-1})$

The sensor model is $P(U_t | R_t)$

*Rain depends on rain in the previous day*
*The rain causes the umbrella to appear*

Note: This network is a *polytree*

# Complete joint distribution

We also need the prior probability distribution at time $0$, $\mathbf{P}(\mathbf{X}_0)$

Putting all together, the **complete joint distribution** over all the variables, for any $t$, computed from the network is:

$$\mathbf{P}(\mathbf{X}_{0:t}, \mathbf{E}_{1:t}) = \mathbf{P}(\mathbf{X}_0) \prod_{i=1}^{t} \mathbf{P}(\mathbf{X}_i|\mathbf{X}_{i-1}) \, \mathbf{P}(\mathbf{E}_i|\mathbf{X}_i)$$

(this formula comes from $P(x_1, \ldots, x_n) = \prod_{i=1}^{n} P(xi \mid parents(Xi))$ )

$\mathbf{P}(\mathbf{X}_0)$        the initial state model over set of variables **X**

$\mathbf{P}(\mathbf{X}_i|\mathbf{X}_{i-1})$     the transition model under a *first-order Markov assumption*

$\mathbf{P}(\mathbf{E}_i|\mathbf{X}_i)$       the sensor model under *Markov sensor assumption*
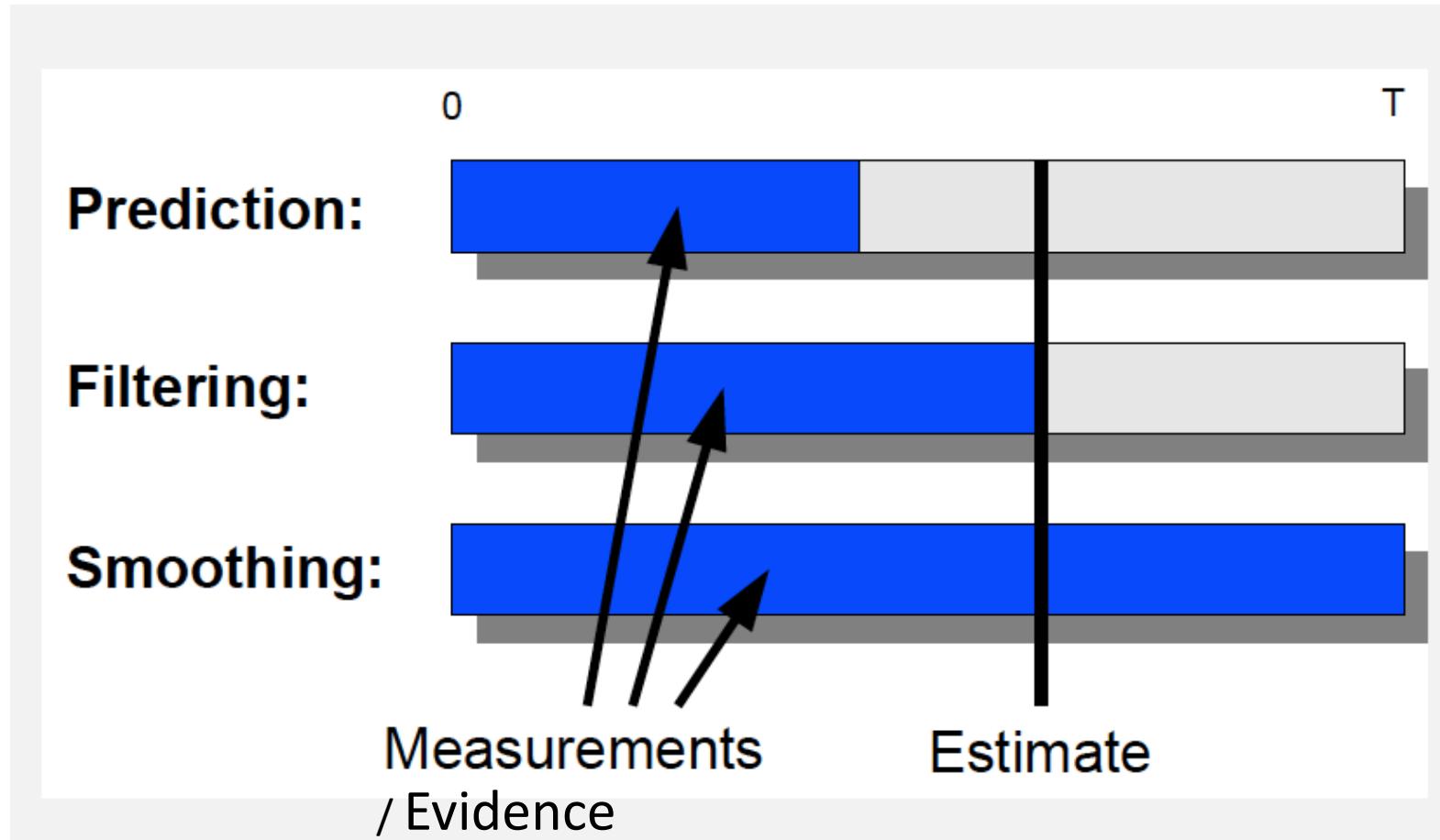
# Inference in temporal models

FILTERING, PREDICTION, SMOOTHING, MOST LIKELY EXPLANATION

# Inference in temporal models

Basic inference tasks based on the temporal model:

- **Filtering**, or **state estimation**: computing the belief state (*posterior probability distribution* of state variables) given evidence from previous and current states. A subtask is the **likelihood of the evidence sequence**. $\mathbf{P}(X_t|e_{1:t})$

- **Prediction:** computing the posterior distribution over a **future** state, given all evidence to date. $\mathbf{P}(X_{t+k}|e_{1:t})$ with $k > 1$.

- **Smoothing:** computing the posterior distribution over a **past** state, given all evidence up to the present. Looking back with the knowledge of today provides a more accurate estimate. $\mathbf{P}(X_k|e_{1:t})$ with $0 \le k < t$.

- **Most likely explanation/sequence**: given a sequence of observations, we might wish to find the sequence of states that is most likely to have generated those observations.

- **Learning**: The transition and sensor models, can be learned from observations.

# Filtering vs other inference tasks



*From a tutorial by Simo Särkkä*

# Filtering: estimation of next state

A good filtering algorithm maintains a **current state estimate** and updates it, rather than going back over the entire history of percepts for each time

The filtering function $f$ takes into account **the state estimation computed up to the present** and the **new evidence.**

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

This process is called **recursive estimation** and is made of two parts:

1. **Prediction**: the current state distribution is projected forward from $t$ to $t+1$:
   $\mathbf{P}(X_{t+1}|e_{1:t})$
2. **Update**: then it is updated using the new evidence $e_{t+1}$: $\mathbf{P}(e_{t+1}|X_{t+1})$

# Computing filtering

It can be computed as follows:

$$\mathbf{P}(X_{t+1}|e_{1:t+1}) = \mathbf{P}(X_{t+1}|e_{1:t}, e_{t+1}) =$$ *Dividing up the evidence*

$$= \alpha\, \mathbf{P}(e_{t+1}|X_{t+1}\, e_{1:t})\, \mathbf{P}(X_{t+1}|e_{1:t})$$ *Bayes's rule, focusing on $e_{t+1}$* ➞

$$= \alpha\, \mathbf{P}(e_{t+1}|X_{t+1})\, \mathbf{P}(X_{t+1}|e_{1:t})$$ *Markov's sensor assumption*

$$= \alpha\, \mathbf{P}(e_{t+1}|X_{t+1}) \sum_{x_t} \mathbf{P}(X_{t+1}|x_t, e_{1:t})\, \mathbf{P}(x_t|e_{1:t})$$ *By conditioning on the current state $X_t$* ➞

$$= \alpha\, \mathbf{P}(e_{t+1}|X_{t+1}) \sum_{x_t} \mathbf{P}(X_{t+1}|x_t)\, \mathbf{P}(x_t|e_{1:t})$$ *Markov's sensor assumption*

Finally: $\mathbf{P}(X_{t+1}|e_{1:t+1}) = \alpha\, \underset{\text{----update----}}{\mathbf{P}(e_{t+1}|X_{t+1})} \underset{\text{---one-step prediction----}}{\sum_{x_t} \mathbf{P}(X_{t+1}|x_t)\, \mathbf{P}(x_t|e_{1:t})}$  (*Filtering*)

We can think of $\mathbf{P}(X_t|e_{1:t})$ as a *message* $f_{1:t}$ that is propagated forward in the sequence. This makes evident the recursive structure:

$f_0 = \mathbf{P}(X_0)$

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

where *Forward* implements the *filtering* update in **costant time**

# Reminder

$$\mathbf{P}(X|Y, \textcolor{red}{Z}) = \frac{\mathbf{P}(Y|X, \textcolor{red}{Z}) \; \mathbf{P}(Y|\textcolor{red}{Z})}{\mathbf{P}(X|\textcolor{red}{Z})} \qquad (\textit{Bayes rule with background knowledge Z})$$

In our case:

$$\mathbf{P}(X_{t+1}|\boldsymbol{e}_{1:t}, \boldsymbol{e}_{t+1}) = \alpha \; \mathbf{P}(\boldsymbol{e}_{t+1}|X_{t+1}\,\boldsymbol{e}_{1:t}) \; \mathbf{P}(X_{t+1}|\boldsymbol{e}_{1:t})$$

**Conditioning** is marginalization applied to conditional probabilities

$$\textit{--- product rule ---}$$

$$\mathbf{P}(\mathbf{Y}) = \sum\nolimits_{Z \in \mathbf{Z}} \mathbf{P}(\mathbf{Y}, \boldsymbol{z}) = \sum\nolimits_{Z \in \mathbf{Z}} \mathbf{P}(\mathbf{Y}|\boldsymbol{z}) \; P(z) \qquad (\textit{Conditioning})$$

In our case:

$$\mathbf{P}(X_{t+1}|\boldsymbol{e}_{1:t}) = \sum\nolimits_{\boldsymbol{x}_t} \mathbf{P}(X_{t+1}|\boldsymbol{x}_t, \boldsymbol{e}_{1:t}) \; \mathbf{P}(\boldsymbol{x}_t|\boldsymbol{e}_{1:t})$$

# Computing filtering for umbrellas

*The example shows the consequences of observing umbrellas for two consecutive days,* $\mathbf{P}(R_2|u_1,u_2)$

**Day 0**: $\mathbf{P}(R_0) = \langle 0.5, 0.5 \rangle$   $f_0$                    no observation, a 50% chance of rain

**Day 1**:

*Predict*: $\mathbf{P}(R_1) = \sum_{r_0} \mathbf{P}(R_1|r_0) \, P(r_0) = \mathbf{P}(R_1|r_0) \times P(r_0) + \mathbf{P}(R_1|\neg r_0) \times P(\neg r_0) =$
$$= \langle 0.7, 0.3 \rangle \times 0.5 + \langle 0.3, 0.7 \rangle \times 0.5 = \langle 0.5, 0.5 \rangle$$

*Update*, observing the evidence $u_1$ ($U_1 = $ true)                                    $f_{1:1}$

$\mathbf{P}(R_1|u_1) = \alpha \, \mathbf{P}(u_1|R_1) \, P(R_1) = \alpha \langle 0.9, 0.2 \rangle \langle 0.5, 0.5 \rangle = \alpha \langle 0.45, 0.1 \rangle \approx \langle 0.818, 0.182 \rangle$

**Day 2**:

*Predict*: $\mathbf{P}(R_2) = \sum_{r_1} \mathbf{P}(R_2|r_1) \, P(r_1) = \langle 0.7, 0.3 \rangle \times 0.818 + \langle 0.3, 0.7 \rangle \times 0.182 = \langle 0.627, 0.373 \rangle$

*Update,* observing the evidence $u_2$                                                         $f_{1:2}$

$\mathbf{P}(R_2|u_1,u_2) = \alpha \, \mathbf{P}(u_2|R_2) \, \mathbf{P}(R_2|u_1) = \alpha \langle 0.9, 0.2 \rangle \langle 0.627, 0.373 \rangle = \alpha \langle 0.565, 0.075 \rangle \approx \langle 0.883, 0.117 \rangle$

The probability of rain increases from day 1 to day 2 because rain persists.

# Prediction

The task of prediction can be seen simply as filtering without the contribution of new evidence. The filtering process already incorporates a **one-step** prediction.

In general, looking ahead $k$ steps, at time $t+k+1$, given evidence up to $t$:

$$\mathbf{P}(\boldsymbol{X}_{t+k+1}|\boldsymbol{e}_{1:t}) = \sum\nolimits_{X_{t+k}} \mathbf{P}(\boldsymbol{X}_{t+k+1}|\boldsymbol{x}_{t+k})\ \mathbf{P}(x_{t+k}|\boldsymbol{e}_{1:t}) \qquad \textit{... looking k steps ahead}$$

This computation involves only *the transition model* and **not** the *sensor model*.

We can show that the predicted distribution for *Rain* converges to a fixed point ⟨0.5, 0.5⟩, after which it remains constant for all time (the **stationary distribution** of the Markov process). The **mixing time** is the time to reach the fixed point.

The more uncertainty there is in the transition model, the shorter will be the *mixing time* and the more the future is obscure.

# Likelihood of evidence sequence (skip)

We can use a forward recursion also to compute the likelihood of an evidence sequence $P(\boldsymbol{e}_{1:t})$, useful if we want to compare two models producing the same evidence sequence.

We can derive a recursive equation similar to filtering and a similar likelihood message

$$\ell_{1:t}(X_t) = P(X_t, e_{1:t})$$

Once we have $\ell_{1:t}(X_t)$ we can compute the likelihood of the evidence sequence by summing out on the values of $X_t$

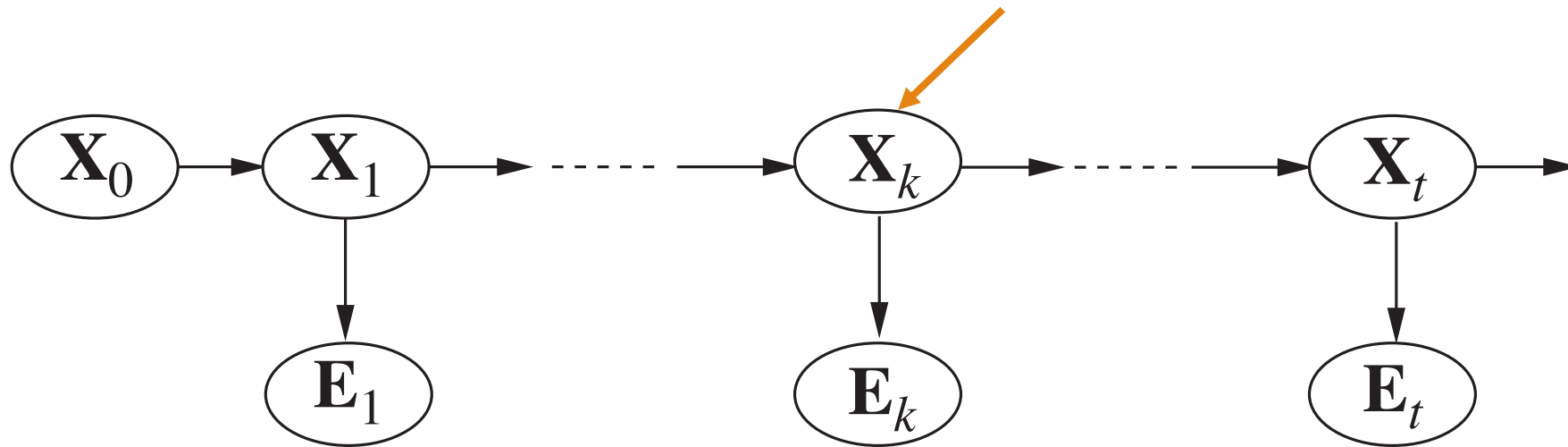$$L_{1:t} = P(\boldsymbol{e}_{1:t}) = \sum_{X_t} \ell_{1:t}(x_t)$$

Note: the likelihood becomes smaller and smaller as $t$ increases.

# Smoothing

**Smoothing** is the process of computing the posterior distribution of the state at **some past time** $k$ given a complete sequence of observations up to the present $t$

$$\mathbf{P}(X_k | e_{1:t}) \qquad \text{for } 0 \le k < t$$

The additional evidence is expected to provide more information and more accurate predictions on the past … The Bayesian network is:

# Computing smoothing

$$\mathbf{P}(X_k | e_{1:t}) = \mathbf{P}(X_k | e_{1:k}, e_{k+1:t}) \qquad \textit{Splitting the evidence}$$

$$= \alpha\, \mathbf{P}(X_k | e_{1:k})\, \mathbf{P}(e_{k+1:t} | X_k, e_{1:k}) \qquad \textit{Bayes's rule}$$

$$= \alpha\, \mathbf{P}(X_k | e_{1:k})\, \mathbf{P}(e_{k+1:t} | X_k) \qquad \textit{Conditional independence}$$

$$= \alpha\, f_{1:k} \times b_{k+1:t} \qquad \times \text{ is pointwise multiplication}$$

$\mathbf{P}(X_k | e_{1:k})$ corresponda to the *forward* message $f_{1:k}$, computed up to $t$, as before in filtering.

We define another *message* $b_{k+1:t} = \mathbf{P}(e_{k+1:t} | X_k)$ that can be computed by a recursive process that runs **backward** from $t$ (next slide).

The terms in $f_{1:k}$ and $b_{k+1:t}$ can be implemented by two recursive calls

- one running forward from 1 to $k$ and using the *filtering equation*

- the other running backward from $t$ to $k+1$ for computing $\mathbf{P}(e_{k+1:t} | X_k)$
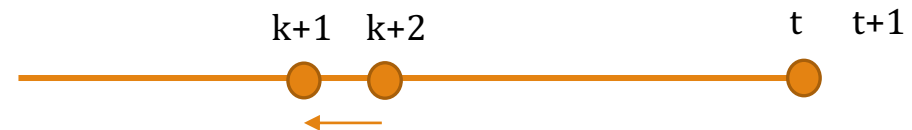
# Computing smoothing going backwards

$\mathbf{P}(\boldsymbol{e}_{k+1:t}|\ X_k) = \sum_{\boldsymbol{x_{k+1}}} \mathbf{P}(\boldsymbol{e}_{k+1:t}|\boldsymbol{X_k}, \boldsymbol{x_{k+1}})\ \mathbf{P}(\boldsymbol{x_{k+1}}|X_k)$ ⬛⬛⬛⬛⬛ *conditioning on* $\boldsymbol{X}_{k+1}$

$\qquad = \sum_{\boldsymbol{x_{k+1}}} \mathbf{P}(\boldsymbol{e}_{k+1:t}|\ x_{k+1})\ \mathbf{P}(\boldsymbol{x_{k+1}}|X_k)$ ⬛⬛⬛ *the evidence only depends on* $\boldsymbol{x}_{k+1}$

$\qquad = \sum_{\boldsymbol{x_{k+1}}} \mathbf{P}(\boldsymbol{e}_{k+1}, \boldsymbol{e}_{k+2:t}|\ x_{k+1})\ \mathbf{P}(\boldsymbol{x_{k+1}}|\boldsymbol{X_k})$ ⬛⬛⬛ *splitting the evidence*

$\qquad = \sum_{\boldsymbol{x_{k+1}}} P(\boldsymbol{e}_{k+1}|\ x_{k+1})\ P(\boldsymbol{e}_{k+2:t}|\ x_{k+1})\ P(\boldsymbol{x_{k+1}}|X_k)$ ⬛⬛ *conditional independence*

$b_{k+1:t} = P(\boldsymbol{e}_{k+2:t}|\ x_{k+1})$ is defined by the previous equation computing backword.

The backward phase starts with $\mathbf{P}(\boldsymbol{e}_{t+1:t}|\ X_t) = \mathbf{1}$ **(a vector of 1's)**, because $(t+1:t)$ is an empty sequence of evidence with probability 1 of observing it.

$\qquad b_{t+1:t} = \mathbf{1}$

$\qquad b_{k+1:t} = Backward(b_{k+2:t}, \boldsymbol{e}_{k+1})$

where *Backward* implements the update defined above

# Computing the smoothed estimate for *rain*

Let's compute the smoothed estimate at time $k=1$ for the probability of rain. Given the observations at time 1 and 2.

$P(R_1 \mid u_1, \underline{u}_2) = \alpha\, P(R_1 \mid u_1) \times P(u_2 \mid R_1)$

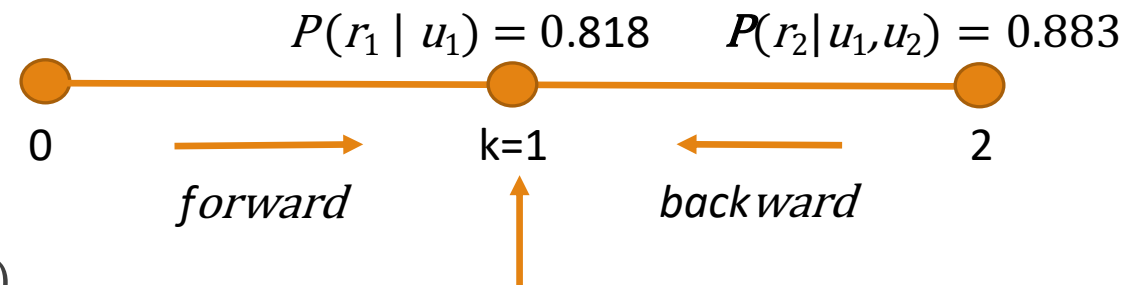$P(R_1 \mid u_1) = \langle 0.818, 0.182 \rangle$

already computed in the filtering stage.

$P(u_2 \mid R_1) = \sum_{r_2} P(u_2 \mid r_2)\, P(U_{3:2} \mid r_2)\, P(r_2 \mid R_1)$

$= (P(u_2 \mid r_2) \times \mathbf{1} \times P(r_2 \mid R_1)) + (P(u_2 \mid \neg r_2) \times \mathbf{1} \times P(\neg r_2 \mid R_1))$

$= (0.9 \times 1 \times \langle 0.7, 0.3 \rangle) + (0.2 \times 1 \times \langle 0.3, 0.7 \rangle) = \langle 0.69, 0.41 \rangle$

$P(R_1 \mid u_1, \underline{u}_2) = \alpha \langle 0.818, 0.182 \rangle \times \langle 0.69, 0.41 \rangle \approx \langle 0.883, 0.117 \rangle$

$P(r_1 \mid u_1) = 0.818 \qquad P(r_2 \mid u_1, u_2) = 0.883$

0 ———→ k=1 ←——— 2

*forward*      *backward*

The **smoothed estimate** for rain on day 1 is higher than the *filtered estimate*. Rain persists.
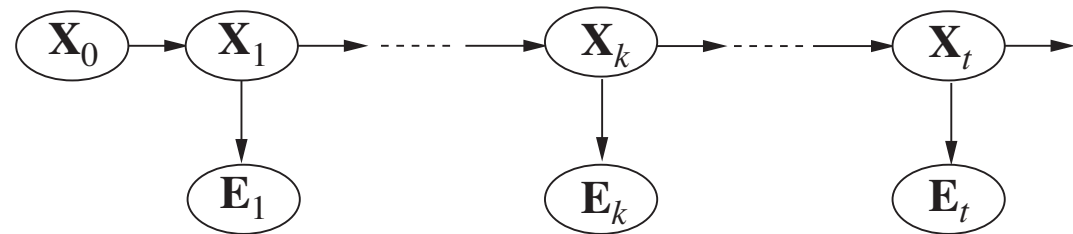
The complexity for smoothing at a specific time $k$ with respect to evidence $e_{1:t}$ is O($t$). For smoothing a sequence O($t^2$). A better linear approach extsts.

# Improving efficiency

The algorithm, called **forward–backward**, uses dynamic programming to reduce the complexity of running the algorithm over the whole sequence.

The trick is to record the results computed during the forward filtering phase, over the whole sequence, and reuse them during the backward phase.

The complexity is O(t) and is consistent with the fact that the Bayesian network structure is a **polytree**.

$$X_0 \rightarrow X_1 \dashrightarrow X_k \dashrightarrow X_t \rightarrow$$
$$X_1 \rightarrow E_1 \qquad X_k \rightarrow E_k \qquad X_t \rightarrow E_t$$

The *forward–backward* algorithm is the basis for many applications that deal with sequences of noisy observations.

Improvements are required to deal with

1. *space complexity* for long sequences
2. *online computations*, where new observations continuously arrive (**fixed-lag smoothing**)

# Forward–backward algorithm for smoothing

**function** FORWARD-BACKWARD($\mathbf{ev}$, $prior$) **returns** a vector of probability distributions
   **inputs**: $\mathbf{ev}$, a vector of evidence values for steps $1, \ldots, t$      *Smoothing all the sequence*
         $prior$, the prior distribution on the initial state, $\mathbf{P}(\mathbf{X}_0)$
   **local variables**: $\mathbf{fv}$, a vector of forward messages for steps $0, \ldots, t$
         $\mathbf{b}$, a representation of the backward message, initially all 1s
         $\mathbf{sv}$, a vector of smoothed estimates for steps $1, \ldots, t$

   $\mathbf{fv}[0] \leftarrow prior$
   **for** $i = 1$ **to** $t$ **do**
      $\mathbf{fv}[i] \leftarrow$ FORWARD($\mathbf{fv}[i-1], \mathbf{ev}[i]$)     *Forward filtering phase*
   **for** $i = t$ **downto** $1$ **do**
      $\mathbf{sv}[i] \leftarrow$ NORMALIZE($\mathbf{fv}[i] \times \mathbf{b}$)     *Backward phase reusing*
      $\mathbf{b} \leftarrow$ BACKWARD($\mathbf{b}, \mathbf{ev}[i]$)        *result from the filtering*
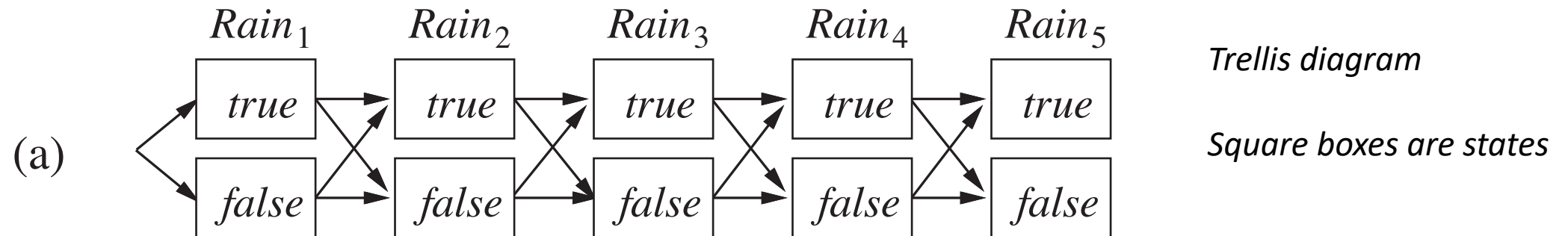   **return** $\mathbf{sv}$

# Finding the most likely sequence

Suppose we observe [*true, true, false, true, true*] for the *Umbrella* variable.

What is the most likely sequence for *Rain* given these oservations?

There are $2^5$ possible *Rain* sequences …Each sequence is a path through a graph whose nodes are the possible states at each time step.

We want to discover which is the one maximizing the likelihood, **in linear time**.

**Likelihood of a path**: product of the transition probabilities along the path and the probabilities of the given observations at each state



(a)

*Trellis diagram*

*Square boxes are states*

# Computing the most likely sequence

There is a recursive relationship between the most likely path to each state $x_{t+1}$ and most likely paths to each previous state $x_t$. We can write a recursive equation, similar to the one for filtering:

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

The forward message in this case is

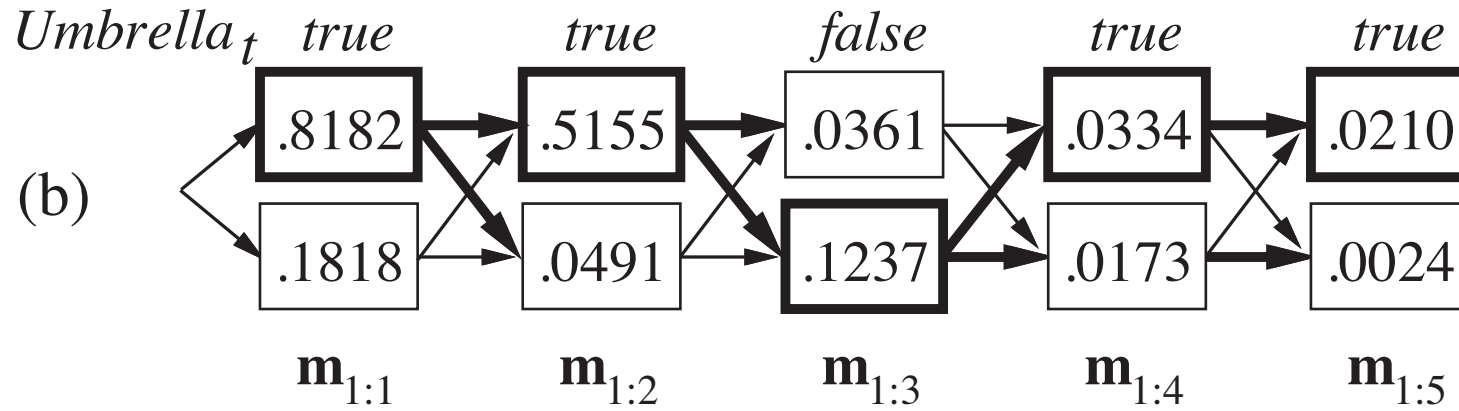$$m_{1:t} = max_{x_1...x_{t-1}} P(x_1, ..., x_{t-1}, X_t \mid e_{1:t})$$

the probabilities of the best path to each state $x_t$; from those we can compute the probabilities of the extended paths at time t+1 and take the max.

The most likely sequence overall can be computed in one pass.

For each state, the best state that leads to it is recorded (marked as black arrows in the example) so that the optimal sequence is identified by following black arrows backwards from the best final state.

This algorithm is the famous **Viterbi algorithm**, named after A. Viterbi [1967]

# Viterbi algorithm on the Umbrella example

$Umbrella_t$   *true*      *true*      *false*      *true*      *true*

(b)

$$\langle .8182 \rangle \rightarrow \langle .5155 \rangle \rightarrow \langle .0361 \rangle \rightarrow \langle .0334 \rangle \rightarrow \langle .0210 \rangle$$

$$\langle .1818 \rangle \rightarrow \langle .0491 \rangle \rightarrow \langle .1237 \rangle \rightarrow \langle .0173 \rangle \rightarrow \langle .0024 \rangle$$

$\mathbf{m}_{1:1}$          $\mathbf{m}_{1:2}$          $\mathbf{m}_{1:3}$          $\mathbf{m}_{1:4}$          $\mathbf{m}_{1:5}$

$\mathbf{m}_{1:1} = \langle 0.8182, 0.1818 \rangle$ as computed in filtering

$\mathbf{m}_{1:2} = \langle \max[P(u_2 \,|\, r_2) \times 0.8182 \times P(r_2 \,|\, r_1), \; P(u_2 \,|\, r_2) \times 0.1818 \times P(r_2 \,|\, \neg r_1)], \quad$ *rain at time 2*

$\qquad \max[P(u_2 \,|\, \neg r_2) \times 0.8182 \times P(\neg r_2 \,|\, r_1), P(u_2 \,|\, \neg r_2) \times 0.1818 \times P(\neg r_2 \,|\, \neg r_1)] \rangle \quad$ *not rain at time 2*

$\qquad = \langle \max[0.9 \times 0.8182 \times 0.7, 0.9 \times 0.1818 \times 0.3], \max[0.2 \times 0.8182 \times 0.3, 0.2 \times 0.1818 \times 0.7] \rangle$

$\qquad = \langle \max[0.5155, 0.0491], \max[0.0491, 0.0255] \rangle = \langle 0.5155, 0.0491 \rangle$

$\mathbf{m}_{1:3} = P(\neg u_3 \,|\, R_3) \times \max_{r2}[P(R_3 \,|\, r_2) \times \mathbf{m}_{1:2}] = \langle 0.361, 0.1237 \rangle$

$\mathbf{m}_{1:4} = P(u_4 \,|\, R_4) \times \max_{r3}[P(R_4 \,|\, r_3) \times \mathbf{m}_{1:3}] = \langle 0.334, 0.173 \rangle$

$\mathbf{m}_{1:5} = P(u_5 \,|\, R_5) \times \max_{r4}[P(R_5 \,|\, r_4) \times \mathbf{m}_{1:4}] = \langle 0.0210, 0.0024 \rangle$

# Applications of the Viterbi algorithm

- The original application of Viterbi was in telecommunications: Viterbi decoders are used for decoding a bitstream encoded using a technique called **convolutional code** or **trellis code**.

- In natural language processing:
    - Different kinds of "sequence tagger": PoS taggers, NE taggers …
    - Dependency parsing
    - Speech recognition (speech-to-text), speech synthesis, speech diarization (recognizing sequences from the same speaker), keyword spotting …
- …

# Other approaches to uncertain reasoning

# An overview of other approaches

Probability theory and Bayesian networks are the dominant approaches today, despite the disadvantage of having to specify many probability values (lots of numbers).

Open issue: Humans reason with numbers? Or their reasoning is more "qualitative"?

Other approaches have been used:

1. Rule-based methods for uncertain reasoning in expert systems
2. Representing ignorance: Dempster–Shafer theory
3. Representing vagueness: fuzzy sets and fuzzy logic

# Rule-based methods for uncertain reasoning

Three good properties of classical logic-based rules:

1. **Locality**: In logical systems, from A and A $\Rightarrow$ B, we can conclude B, without worrying about any other rules. In probabilistic systems, we need to consider all the evidence.

2. **Detachment**: Once B is proved, it can be used regardless of how it was derived. It can be *detached* from its justification. Probabilities cannot be detached from evidence.

3. **Truth-functionality**: the truth of complex sentences can be computed from the truth of the components. Probability combination does not work this way.

   Ex  $P(H_1) = 0.5$, $P(H_2) = 0.5$ but $P(H_1 \vee H_2) = 0.75$

The idea is to attach **degree of belief** to facts and rules and to combine and propagate them

The most famous example is the **certainty factors model**, which was developed for the MYCIN medical diagnosis program. The system was carefully engineered to avoid mixing different kind of rules (diagnostic vs causal), trying to control non plausible results.
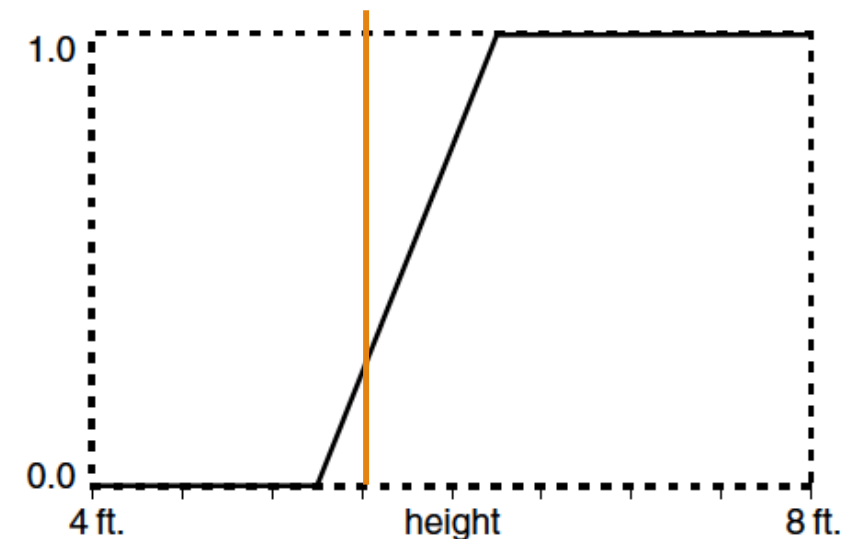
Part of the mechanism was later explained in terms of a limited form of probabilistic reasoning.

# Fuzzy set theory: representing vagueness

The theory of **fuzzy sets** is a way to specify **how well** an object satisfies a vague description (non categorical) property, like "being tall".

✓ Crispy sets: Nate is tall is True or False. Membership function: $\boldsymbol{\mu}$(S) in {0, 1}

✓ Fuzzy sets: Nate is 1.75m. Is Nate tall? Sort of ... the are degrees of tallness.
  - Membership function: $\boldsymbol{\mu}$(S) in [0, 1].
  - 0 means definetely not tall; 1 definetely tall.
  - The set does not have sharp boundaries.

Fuzziness id not uncertainty in the world, uncertainty in the use of qualifiers/properties.



Relation between degree of 'tallness' and height

# Fuzzy logic

Fuzzy logic is way of reasoning about membership in fuzzy sets.

A fuzzy predicate implicitly defines a **fuzzy set**.

Fuzzy logic is a method for reasoning with logical expressions describing membership in fuzzy sets.

Standard rules:

$\text{T}(A \wedge B) = min(\text{T}(A), \text{T}(B))$     *A, B are assertions or properties*

$\text{T}(A \vee B) = max(\text{T}(A), \text{T}(B))$     *they denote degree of membership in a set*

$\text{T}(\neg A) = 1 - \text{T}(A)$

Example. $\text{T}(Tall(john) \wedge Heavy(john)) = \text{MIN}(\text{T}(Tall(john), \text{T}(Heavy(john)))$

# Fuzzy control

Fuzzy control is a methodology for constructing control systems in which the mapping between real-valued input and output parameters is represented by fuzzy rules.
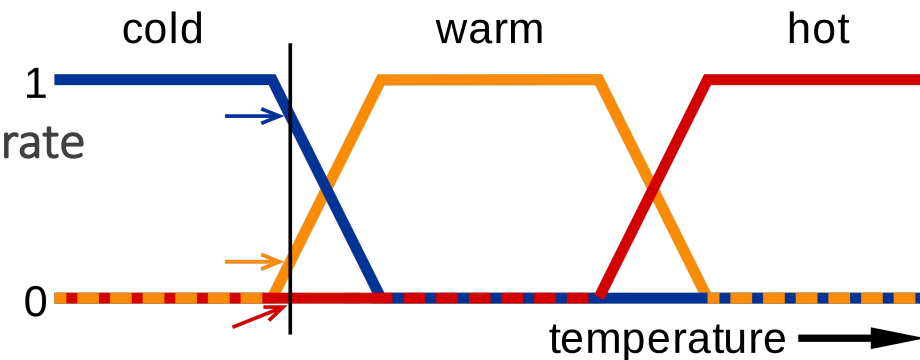
1. **Fuzzification**: continous variables are mapped into a set of linguistic variables with fuzzy values; *temperature* can be mapped in fuzzy variables such as *low, medium, high* and their membership functions

2. **Reasoning** with rules expressed in terms of these fuzzy variables, reaching fuzzy conclusions.
    IF temperature IS cold THEN fan_speed is slow
    IF temperature IS warm THEN fan_speed is moderate
    IF temperature IS hot THEN fan_speed is high

1. **Defuzzification**: map the results into numerical values

cold          warm          hot

1

0

temperature →

# Conclusions

- A dynamic world can be represented using a set of random variables to represent the state at each point in time. A special case of Bayesian network.

- A temporal probability model can be thought of as containing a **transition model** describing the state evolution and a **sensor model** describing the observation process.

- Assuming the **Markov property** and **stationary processes** greatly simplifies the representation and the computation.

- The principal inference tasks in temporal models are: **filtering**, **prediction**, and **computing the most likely explanation sequence**. Simple, recursive algorithms exist *linear* in the length of the sequence.

- Algorithms can be further simplified assuming **Hidden Markov Models**, i.e. states are represented by a single discrete variable.

- **Hidden Markov models**, **Kalman filters**, and **dynamic Bayesian networks** are left out.

# References

- ✓ Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach* (3<sup>rd</sup> edition). Pearson Education 2010 [Cap 15 – Probabilistic reasoning over time]