

Probabilistic reasoning over time -

# Hidden Markov Models

(recap BNs)

Applied artificial intelligence (EDA132)

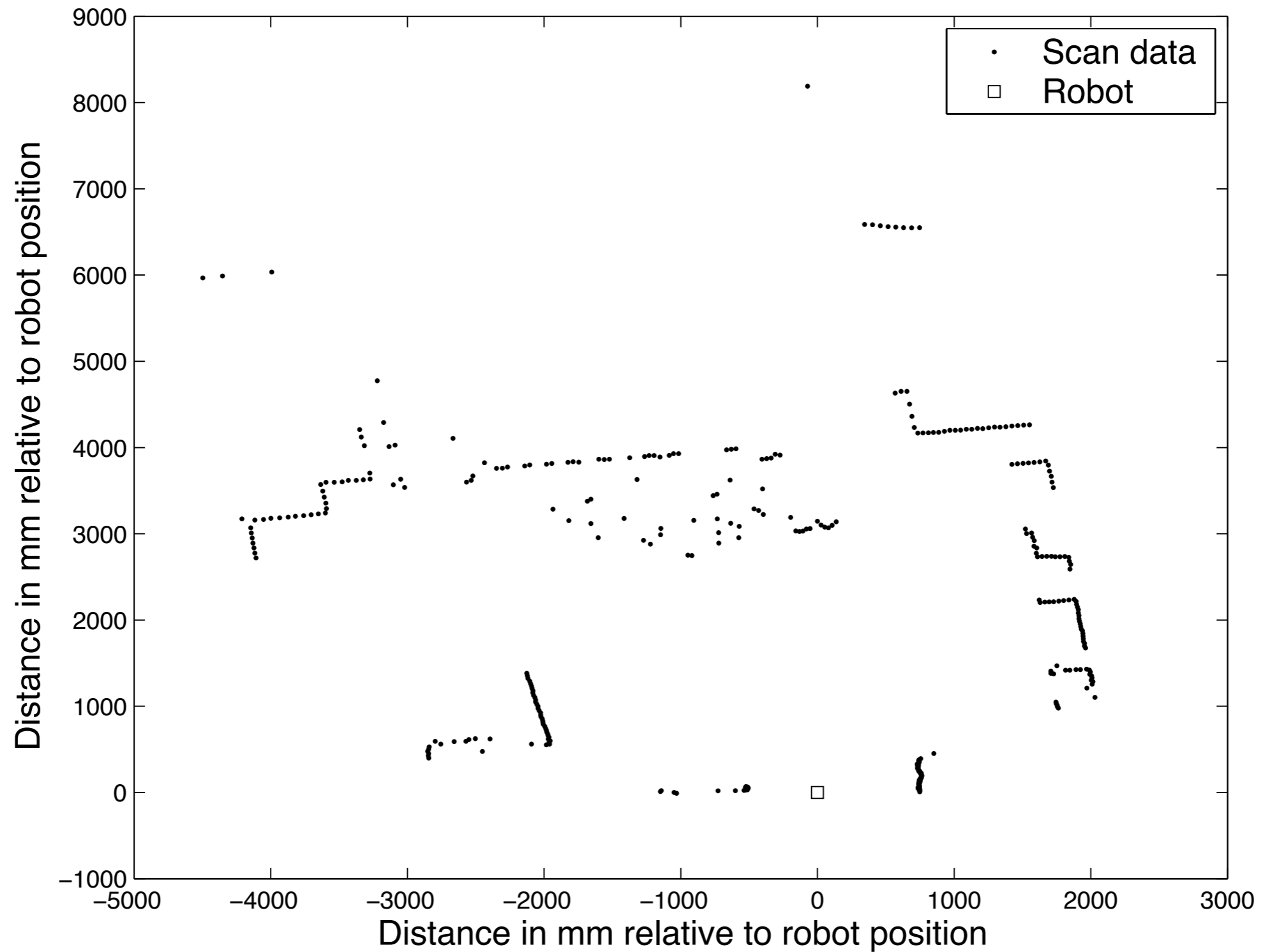
Lecture 10

2016-02-17

Elin A. Topp

Material based on course book, chapter 15

# A robot's view of the world...

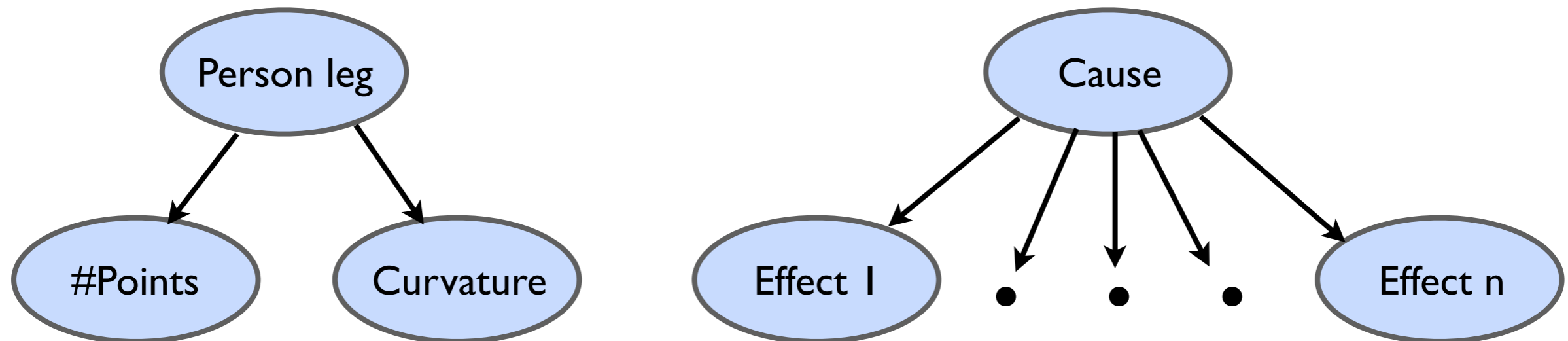


# Bayes' Rule and conditional independence

$$\begin{aligned} & \mathbb{P}( \text{PersonLeg} \mid \#pointsInRange \wedge \text{curvatureCorrect} ) \\ &= \alpha \mathbb{P}( \#pointsInRange \wedge \text{curvatureCorrect} \mid \text{PersonLeg} ) \mathbb{P}( \text{PersonLeg} ) \\ &= \alpha \mathbb{P}( \#pointsInRange \mid \text{PersonLeg} ) \mathbb{P}( \text{curvatureCorrect} \mid \text{PersonLeg} ) \mathbb{P}( \text{PersonLeg} ) \end{aligned}$$

An example of a *naive Bayes* model:

$$\mathbb{P}( \text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n ) = \mathbb{P}( \text{Cause} ) \prod_i \mathbb{P}( \text{Effect}_i \mid \text{Cause} )$$



The total number of parameters is *linear* in  $n$

# Bayesian networks

A simple, graphical notation for *conditional independence assertions* and hence for compact specification of full joint distributions

Syntax:

- a set of nodes, one per random variable

- a directed, acyclic graph (link  $\approx$  “directly influences”)

- a conditional distribution for each node given its parents:

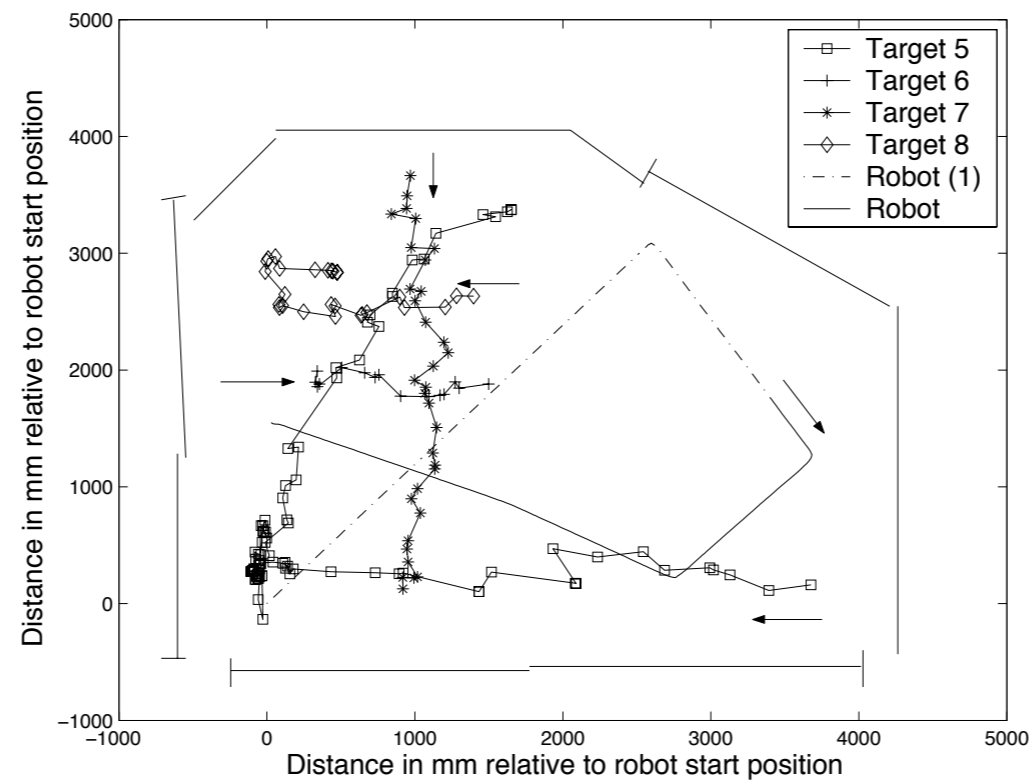
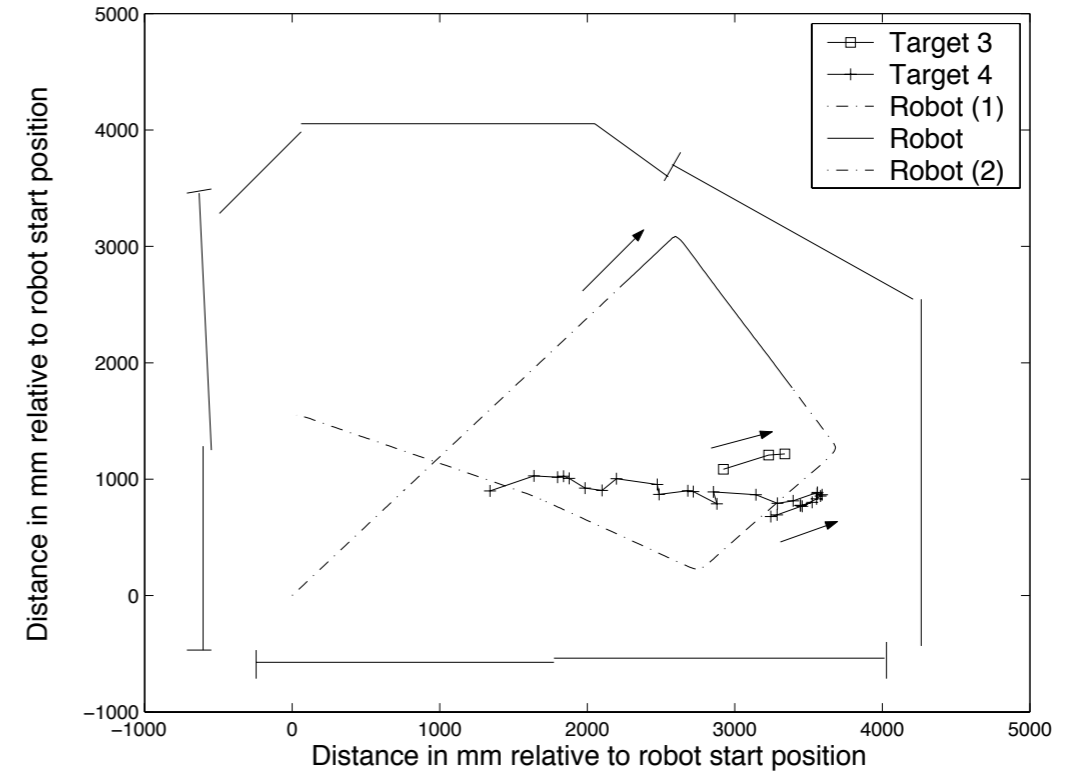
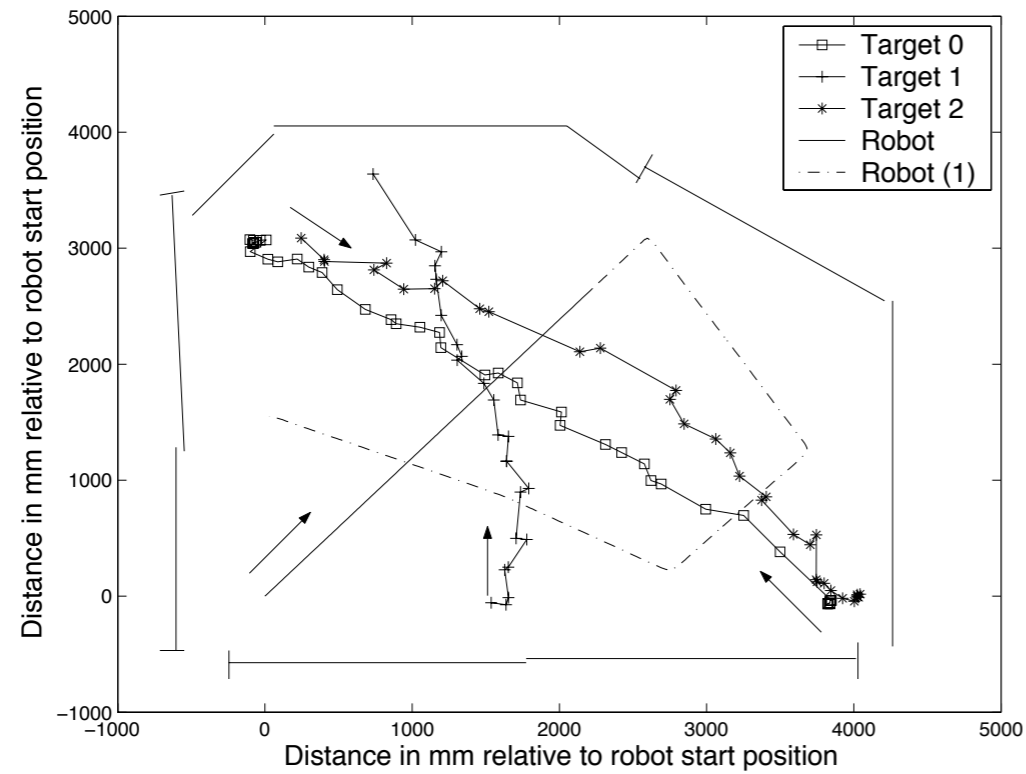
$$P(X_i | Parents(X_i))$$

In the simplest case, conditional distribution represented as a

*conditional probability table* ( CPT)

giving the distribution over  $X_i$  for each combination of parent values

# Tracking and associating... while moving ...



# Probabilistic reasoning over time

... means to keep track of the current state of

- a process (temperature controller, other controllers)
- an agent with respect to the world (localisation of a robot in some “world”)

in order to make predictions or to simply understand what might have caused this current state.

This involves both a **transition model** (how the state is assumed to change) and a **sensor model** (how observations / percepts are related to the world state).

Previously:

the focus was on what was possible to happen (e.g., search), now it is on what is likely / unlikely to happen

the focus was on static worlds (Bayesian networks), now we look at dynamic processes where everything (state AND observations) depend on time.

# Three classes of approaches

Hidden Markov models

(Particle filters)

Kalman filters

Dynamic Bayesian networks (cover actually the other two as special cases)

But first, some basics ...

# Reasoning over time

With

$X_t$  the current state description at time  $t$

$E_t$  the evidence obtained at time  $t$

we can describe a *state transition model* and a *sensor model* that we can use to model a time step sequence - a chain of states and sensor readings according to discrete time steps - so that we can understand the ongoing process.

We assume to start out in  $X_0$ , but evidence will only arrive after the first state transition is made:  $E_1$  is then the first piece of evidence to be plugged into the chain.

The “general” transition model would then specify

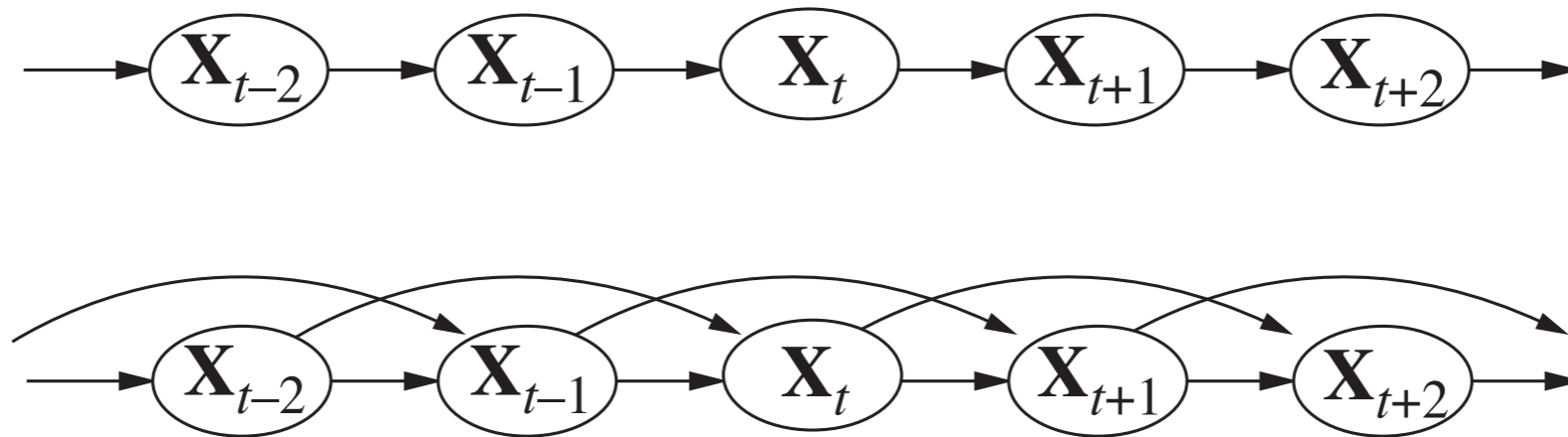
$$\mathbb{P}(X_t | X_{0:t-1})$$

... this would mean we need full joint distributions over all time steps... or not?

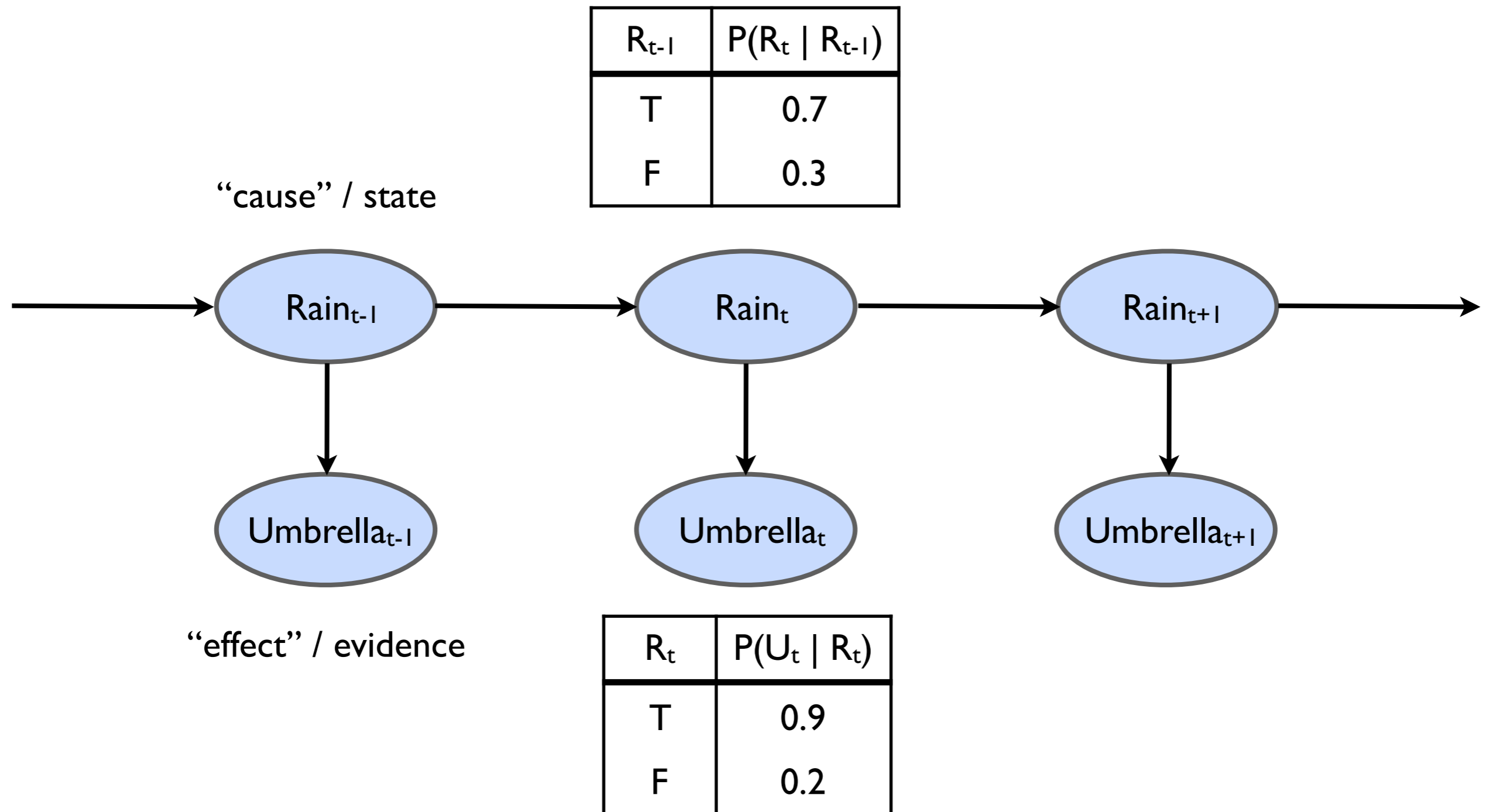


# The Markov assumption

A process is *Markov* (i.e., complies with the Markov assumption), when any given state  $X_t$  depends only on a *finite and fixed number of previous states*.



# A first-order Markov chain as Bayesian network



# Inference for any t

With

$\mathbb{P}(X_0)$  the prior probability distribution in  $t=0$  (i.e., the *initial state model*),

$\mathbb{P}(X_i | X_{i-1})$  the state transition model and

$\mathbb{P}(E_i | X_i)$  the sensor model

we have the complete joint distribution for all variables for any t.

$$\mathbb{P}(X_{0:t}, E_{1:t}) = \mathbb{P}(X_0) \prod_{i=1}^t \mathbb{P}(X_i | X_{i-1}) \mathbb{P}(E_i | X_i)$$

# The Markov assumption

First-order Markov chain:

State variables (at  $t$ ) contain ALL information needed for  $t+1$ .

Sometimes, that is too strong an assumption (or too weak in some sense).

Hence, increase either the order (second-order Markov chain)

or

add information into the state variable(s) ( $R$  could include also *Season, Humidity, Pressure, Location*, instead of only “*Rain*”)

Note: It is possible to express an increase in order by increasing the number of state variables, keeping the order fixed - for the umbrella world you could use

$R = \langle \text{RainYesterday}, \text{RainToday} \rangle$

When things get too complex, rather add another sensor (e.g., observe coats).

# Inference in temporal models

## - what can we use all this for?

- **Filtering:** Finding the **belief state**, or doing **state estimation**, i.e., computing the posterior distribution over the *most recent state*, using evidence up to this point:  $\mathbb{P}(X_t | e_{1:t})$
- **Predicting:** Computing the posterior over a *future state*, using evidence up to this point:  $\mathbb{P}(X_{t+k} | e_{1:t})$  for some  $k > 0$  (can be used to evaluate course of action based on predicted outcome)
- **Smoothing:** Computing the posterior over a past state, i.e., understand the past, given information up to this point:  $\mathbb{P}(X_k | e_{1:t})$  for some  $k$  with  $0 \leq k < t$
- **Explaining:** Find the best explanation for a series of observations, i.e., computing  $\operatorname{argmax}_{x_{1:t}} P(x_{1:t} | e_{1:t})$  - can be efficiently handled by Viterbi algorithm
- **Learning:** If sensor and / or transition model are not known, they can be learned from observations (by-product of inference in Bayesian network - both static or dynamic). Inference gives estimates, estimates are used to update the model, updated models provide new estimates (by inference). Iterate until converging - again, this is an instance of the EM-algorithm.

# Filtering: Prediction & update (FORWARD-step)

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

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

(decompose)

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

(Bayes' Rule)

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

(1. Markov assumption (sensor model),  
2. one-step prediction)

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

(sum over atomic events for  $X$ )

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

(Markov assumption)

$$\mathbb{P}(X_t | e_{1:t})$$

(“forward message”, propagated recursively)

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

through “forward step function”)

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

# Prediction - filtering without the update

$$\mathbb{P}(\mathbf{X}_{t+k+1} | \mathbf{e}_{1:t}) = \sum_{\mathbf{x}_{t+k}} \mathbb{P}(\mathbf{X}_{t+k+1} | \mathbf{x}_t) P(\mathbf{x}_{t+k} | \mathbf{e}_{1:t}) \quad (\text{k-step prediction})$$

For large  $k$  the prediction gets quite blurry and will eventually converge into a *stationary distribution* at the *mixing point*, i.e., the point in time when this convergence is reached - in some sense this is when “everything is possible”.

# Smoothing: “explaining” backward

$\mathbb{P}(X_k | e_{1:t}) = fb(X_k, e_{1:k}, \mathbb{P}(e_{k+1:t} | X_k))$  with  $0 \leq k < t$  (understand the past from the recent past)

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

$= \alpha \mathbb{P}(X_k | e_{1:k}) \mathbb{P}(e_{k+1:t} | X_k, e_{1:k})$  (Bayes' Rule)

$= \alpha \mathbb{P}(X_k | e_{1:k}) \mathbb{P}(e_{k+1:t} | X_k)$  (Markov assumption)

$= \alpha f_{1:k} \times b_{k+1:t}$  (forward-message  $\times$  backward-message)



# Smoothing: calculating backward message

$$b_{k+1:t} = \mathbb{P}(e_{k+1:t} | \mathbf{X}_k)$$

$$= \sum_{\mathbf{x}_{k+1}} \mathbb{P}(e_{k+1:t} | \mathbf{X}_k, \mathbf{x}_{k+1}) \mathbb{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \quad (\text{conditioning on } \mathbf{X}_{k+1}, \text{ i.e., looking "backward"})$$

$$= \sum_{\mathbf{x}_{k+1}} P(e_{k+1:t} | \mathbf{x}_{k+1}) \mathbb{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \quad (\text{cond. indep. - Markov assumption})$$

$$= \sum_{\mathbf{x}_{k+1}} P(e_{k+1}, e_{k+2:t} | \mathbf{x}_{k+1}) \mathbb{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \quad (\text{decompose})$$

$$= \sum_{\mathbf{x}_{k+1}} P(e_{k+1} | \mathbf{x}_{k+1}) P(e_{k+2:t} | \mathbf{x}_{k+1}) \mathbb{P}(\mathbf{x}_{k+1} | \mathbf{X}_k) \quad (1. \text{ sensor, 2. backward msg, 3. transition model})$$

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

$$\mathbb{P}(e_{k+1:t} | \mathbf{X}_k) \quad (\text{"backward message", propagated recursively})$$

$$b_{k+1:t} = \text{BACKWARD}(b_{k+2:t}, e_{k+1}) \quad (\text{through "backward step function"})$$

$$b_{t+1:t} = \mathbb{P}(e_{t+1:t} | \mathbf{X}_t) = \mathbb{P}(\cdot | \mathbf{X}_t) = I$$

# Smoothing “in a nutshell”: Forward-Backward-algorithm

$\mathbb{P}(\mathbf{X}_k | \mathbf{e}_{1:t}) = fb(\mathbf{e}_{1:k}, \mathbb{P}(\mathbf{e}_{k+1:t} | \mathbf{X}_k))$  with  $0 \leq k < t$     understand the past from the recent past

$$= \alpha \mathbf{f}_{1:k} \times \mathbf{b}_{k+1:t}$$

by first filtering (forward) until step  $k$ , then explaining backward from  $t$  to  $k+1$

Obviously, it is a good idea to store the filtering (forward) results for later smoothing

Drawback of the algorithm: not really suitable for online use ( $t$  is growing, ...)

Consequently, try with fixed-lag-smoothing (keeping a fixed-length window, BUT: “simple”  
Forward-Backward does not really do it efficiently - here we need HMMs)

# “HMM”

## Hidden Markov models

A specific class of models (sensor and transition) to be plugged into the previously discussed algorithms - which makes the algorithms more specific as well!

### Main idea:

The state is represented by a *single discrete random variable*, taking on values that represent the (all) possible states of the world.

Complex states, e.g., the location and the heading of a robot in a grid world can be merged into one variable; the possible values are then all possible tuples of the values for each original “single” variable.

# “HMM”

## State transition and sensor model

We get the following notation:

$X_t$  the state at time  $t$ , taking on values  $1 \dots S$ , with  $S$  the number of possible states / values.

$E_t$  the observation at time  $t$

The *transition* model  $P(X_t | X_{t-1})$  is then expressed as  $S \times S$  matrix  $T$ :

$$T_{ij} = P(X_t = j | X_{t-1} = i) \text{ in time step } t$$

The *sensor* model for the corresponding observations depending on the current state, i.e.,  $P(e_t | X_t = i)$  is then expressed as  $S \times S$  diagonal matrix  $O$  in time step  $t$  with

$$O_{e_t, ij}^{e_t} = P(e_t | X_t = i) \quad \text{for } i = j \quad \text{and}$$

$$O_{e_t, ij}^{e_t} = 0 \quad \text{for } i \neq j$$

# Forward-backward equations as matrix-vector operations

Forward-equation (recap)

$$P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) = f(P(\mathbf{X}_t | \mathbf{e}_{1:t}), \mathbf{e}_{t+1}) = \mathbf{f}_{1:t+1} = \alpha P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t})$$

becomes  $\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^T \mathbf{f}_{1:t}$

Backward-equation (recap)

$$P(\mathbf{e}_{k+1:t} | \mathbf{X}_k) = \mathbf{b}_{k+1:t} = \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1} | \mathbf{x}_{k+1}) P(\mathbf{e}_{k+2:t} | \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} | \mathbf{X}_k)$$

becomes  $\mathbf{b}_{k+1:t} = \mathbf{T} \mathbf{O}_{k+1} \mathbf{b}_{k+2:t}$

Forward-Backward-equation is then still  $\propto \mathbf{f}_{1:k} \times \mathbf{b}_{k+1:t}$

# Smoothing in constant space

## Idea

propagate both  $f$  and  $b$  in the same direction, hence avoiding to store the  $f_{l:k}$  for a shifting / growing time slice  $k:t$

Propagate the forward-message  $f$  “backward” with

$$f_{l:t} = \alpha' (T^T)^{-1} O^{-1}_{t+1} f_{l:t+1}$$

Start with computing  $f_{t:t}$  in a standard forward-run, forgetting all the intermediate messages, then compute both  $f$  and  $b$  simultaneously “backward” to do smoothing for each step this should be done for (NOTE: works obviously only if  $T^T$  and  $O$  can be inverted, i.e., every sensor reading must be possible in every state, though it can be very unlikely)

# Fixed-lag smoothing (online)

## Idea

if we can do smoothing with constant space requirements, we can also find an efficient recursive algorithm for online smoothing (a shifting “window”), independent of the length  $d$  of the investigated time slice  $t-d$  (with  $t$  growing).

We need to compute

$\propto f_{l:t-d} \times b_{t-d+1:t}$  for time slice  $t-d$ . In  $t+1$ , when a new observation arrives, we need

$\propto f_{l:t-d+1} \times b_{t-d+1:t+1}$  for time slice  $t-d+1$ .

We can get  $f_{l:t-d+1}$  from  $f_{l:t-d}$ , applying standard filtering.

For the backward message, some more inspection has to be done ( $b_{t-d+1:t+1}$  depends on the new evidence in  $t+1$ ) but there is a way by looking at how  $b_{t-d+1:t}$  relates to  $b_{t+1:t}$

# Fixed-lag smoothing (online)

Backward recursion:

apply the recursive equation for  $\mathbf{b}_{t-d+1:t}$   $d$  times:

$$\mathbf{b}_{t-d+1:t} = \left( \prod_{i=t-d+1}^t \mathbf{TO}_i \right) \mathbf{b}_{t+1:t} = \mathbf{B}_{t-d+1:t} \mathbf{l}$$

Then, after the next observation, this will be:

$$\mathbf{b}_{t-d+2:t+1} = \left( \prod_{i=t-d+2}^{t+1} \mathbf{TO}_i \right) \mathbf{b}_{t+2:t+1} = \mathbf{B}_{t-d+2:t+1} \mathbf{l}$$

Do some matrix “division” and get an incremental update for  $\mathbf{B}$  (and ultimately  $\mathbf{b}_{t-d+2:t+1}$ ):

$$\mathbf{B}_{t-d+2:t+1} = \mathbf{O}^{-1}_{t-d+1} \mathbf{T}^{-1} \mathbf{B}_{t-d+1:t} \mathbf{TO}_{t+1}$$



# The full algorithm for fixed-lag smoothing

**function** FIXED-LAG-SMOOTHING( $e_t, hmm, d$ ) **returns** a distribution over  $\mathbf{X}_{t-d}$

**inputs:**  $e_t$ , the current evidence for time step  $t$

$hmm$ , a hidden Markov model with  $S \times S$  transition matrix  $\mathbf{T}$

$d$ , the length of the lag for smoothing

**persistent:**  $t$ , the current time, initially 1

$\mathbf{f}$ , the forward message  $\mathbf{P}(X_t|e_{1:t})$ , initially  $hmm.PRIOR$

$\mathbf{B}$ , the  $d$ -step backward transformation matrix, initially the identity matrix

$e_{t-d:t}$ , double-ended list of evidence from  $t-d$  to  $t$ , initially empty

**local variables:**  $\mathbf{O}_{t-d}, \mathbf{O}_t$ , diagonal matrices containing the sensor model information

add  $e_t$  to the end of  $e_{t-d:t}$

$\mathbf{O}_t \leftarrow$  diagonal matrix containing  $\mathbf{P}(e_t|X_t)$

**if**  $t > d$  **then**

$\mathbf{f} \leftarrow$  FORWARD( $\mathbf{f}, e_t$ )

remove  $e_{t-d-1}$  from the beginning of  $e_{t-d:t}$

$\mathbf{O}_{t-d} \leftarrow$  diagonal matrix containing  $\mathbf{P}(e_{t-d}|X_{t-d})$

$\mathbf{B} \leftarrow \mathbf{O}_{t-d}^{-1} \mathbf{T}^{-1} \mathbf{B} \mathbf{O}_t$

**else**  $\mathbf{B} \leftarrow \mathbf{B} \mathbf{O}_t$

$t \leftarrow t + 1$

**if**  $t > d$  **then return** NORMALIZE( $\mathbf{f} \times \mathbf{B}1$ ) **else return** null

# Summary

## Inference in temporal models

- Filtering and prediction (FORWARD)
- Smoothing (FORWARD-BACKWARD)

## Hidden Markov Models

- Simplified matrix representation for Forward-backward calculations

# Assignment 3



(look on the course page...)