

# Probabilistic Graphical Models

## Lecture 7: Dynamic State Space Models

Volkan Cevher, Matthias Seeger  
Ecole Polytechnique Fédérale de Lausanne

21/10/2011

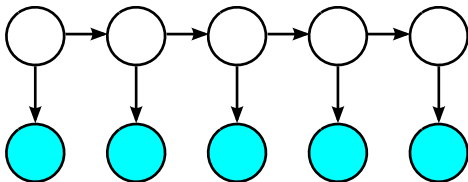


# Announcement

Please bring your *Assignment 5* sheet along to next tutorial.  
Points were not recorded.

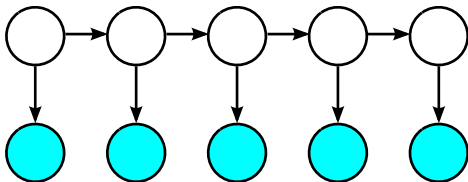
- 1 Hidden Markov Models
- 2 Linear Dynamical Systems
- 3 General Filtering / Smoothing

# Forward in Time



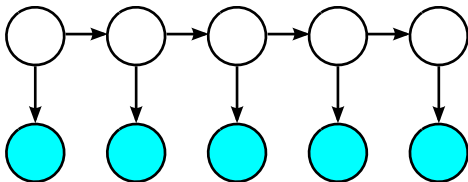
- Final lecture in part I: We'll do a big step ... in time

# Forward in Time



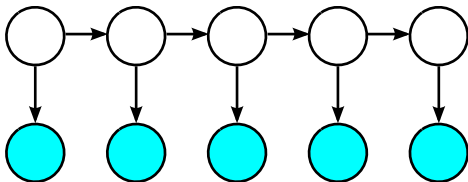
- Final lecture in part I: We'll do a big step . . . in time
- Why dynamic models?
  - Number one reason for causal dependence: Succession in time
  - Filtering, tracking, forward prediction, time series, sequential learning, . . .

# Forward in Time



- Final lecture in part I: We'll do a big step . . . in time
- Why dynamic models?
  - Number one reason for causal dependence: Succession in time
  - Filtering, tracking, forward prediction, time series, sequential learning, . . .
- What's special about dynamic models?
  - Only **one** direction (time arrow)  $\rightarrow$  Linear (in)dependence  $\rightarrow$  Markov chain  $\rightarrow$  Chains are (simple) trees  $\rightarrow$  Belief propagation!

# Forward in Time



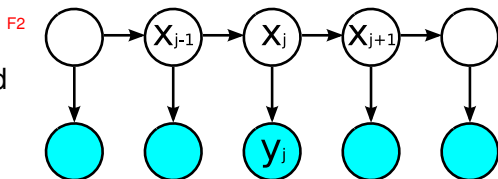
- Final lecture in part I: We'll do a big step . . . in time
- Why dynamic models?
  - Number one reason for causal dependence: Succession in time
  - Filtering, tracking, forward prediction, time series, sequential learning, . . .
- What's special about dynamic models?
  - Only **one** direction (time arrow)  $\rightarrow$  Linear (in)dependence  $\rightarrow$  Markov chain  $\rightarrow$  Chains are (simple) trees  $\rightarrow$  Belief propagation!
- **Markov chain**: Present separates between past and future

$$(\mathbf{x}_{<i}) \perp (\mathbf{x}_{>i}) \mid \mathbf{x}_i$$

# Vocabulary

## Hidden Markov Model

- $P(\mathbf{y}_j | \mathbf{x}_j)$  Observation likelihood  
 $P(\mathbf{x}_j | \mathbf{x}_{j-1})$  Transition kernel  
 $P(\mathbf{x}_1)$  Initial state prior



Stationary Model: CPTs independent of  $j$ .

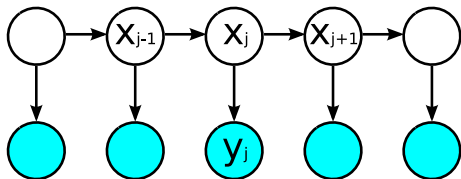
Notation:  $\mathbf{x}_{<j} = (\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$



# Vocabulary

## Hidden Markov Model

- $P(\mathbf{y}_j|\mathbf{x}_j)$  Observation likelihood  
 $P(\mathbf{x}_j|\mathbf{x}_{j-1})$  Transition kernel  
 $P(\mathbf{x}_1)$  Initial state prior



Stationary Model: CPTs independent of  $j$ .

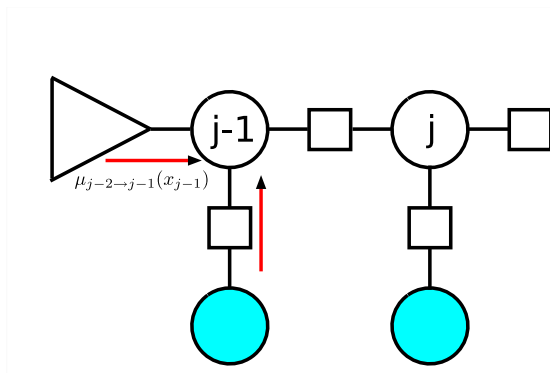
Notation:  $\mathbf{x}_{<j} = (\mathbf{x}_1, \dots, \mathbf{x}_{j-1})$

- Filtering:  $P(\mathbf{x}_j|\mathbf{y}_{\leq j})$  (sequential prediction)
- Smoothing:  $P(\mathbf{x}_j|\mathbf{y}_{1\dots J})$  (inference given past and future)
- Learning: Fitting parameters of  $P(\mathbf{y}|\mathbf{x})$ ,  $P(\mathbf{x}|\mathbf{x}_{-1})$ :  
EM, based on smoothing [EM comes from HMM research]

# Filtering: Information Forward Propagation

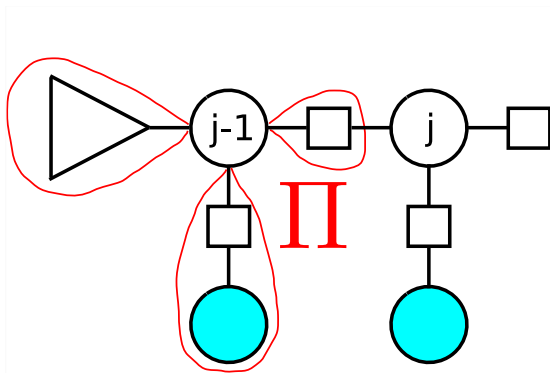
- **Formal:** Message passing on factor graph

F3



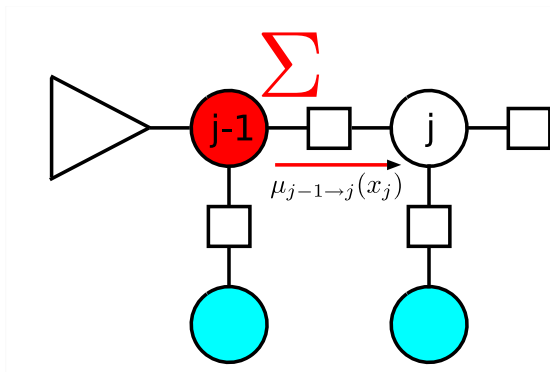
# Filtering: Information Forward Propagation

- **Formal:** Message passing on factor graph



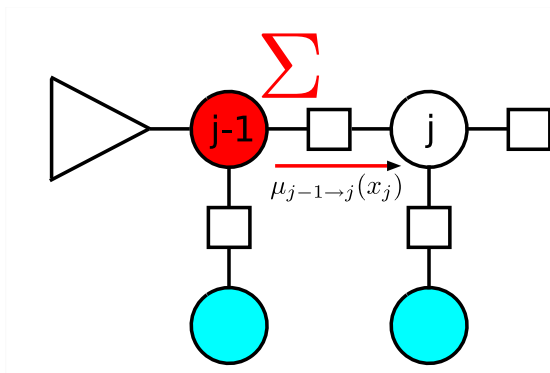
# Filtering: Information Forward Propagation

- **Formal:** Message passing on factor graph



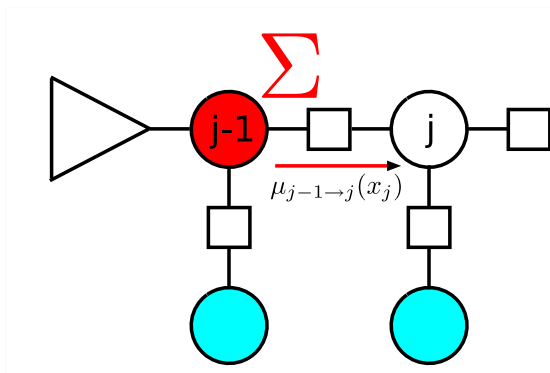
# Filtering: Information Forward Propagation

- **Formal:** Message passing on factor graph
- **Above formulae:** Information forward propagation:  
 $\mu_{(j-2) \rightarrow (j-1)}(x_{j-1})$ :  
 Prior "from the past"



# Filtering: Information Forward Propagation

- **Formal:** Message passing on factor graph
- **Above formulae:** Information forward propagation:  
 $\mu_{(j-2) \rightarrow (j-1)}(x_{j-1})$ :  
 Prior "from the past"



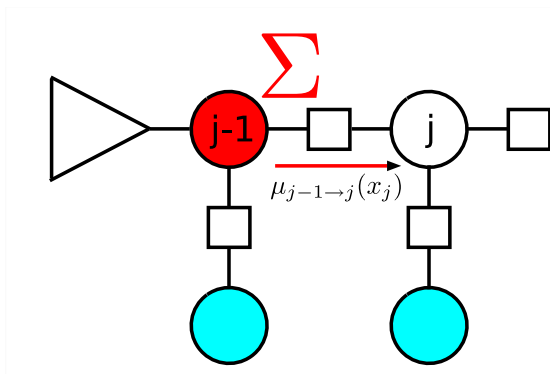
- 1 Measurement: Prior  $\rightarrow$  posterior:  $\propto \mu_{(j-2) \rightarrow (j-1)}(x_{j-1})P(y_{j-1}|x_{j-1})$

# Filtering: Information Forward Propagation

- **Formal:** Message passing on factor graph

- **Above formulae:** Information forward propagation:

$\mu_{(j-2) \rightarrow (j-1)}(x_{j-1})$ :  
Prior "from the past"



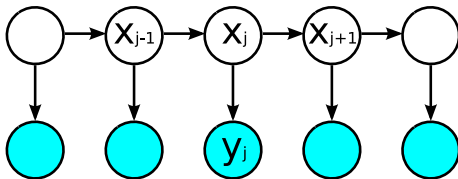
- 1 Measurement: Prior  $\rightarrow$  posterior:  $\propto \mu_{(j-2) \rightarrow (j-1)}(x_{j-1})P(y_{j-1}|x_{j-1})$
- 2 Diffusion, information propagation (marginalization):

$$\mu_{(j-1) \rightarrow j}(x_j) \propto \sum_{x_{j-1}} (\mu_{(j-2) \rightarrow (j-1)}(x_{j-1})P(y_{j-1}|x_{j-1}))P(x_j|x_{j-1})$$

$\Rightarrow$  Know how to do (1), (2)? Can do BP!

F3b

## Sum-Product Algorithm for HMMs



- Backward messages: Exactly the same, just reverse time arrow

① Measurement: Prior  $\rightarrow$  posterior:  $\propto \mu_{j \leftarrow (j+1)}(x_j) P(y_j | x_j)$

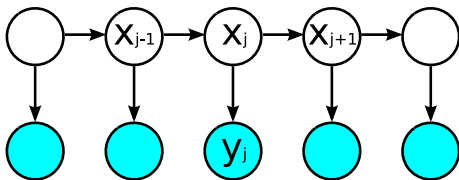
② Diffusion, information propagation (marginalization):

$$\mu_{(j-1) \leftarrow j}(x_{j-1}) \propto \sum_{x_j} P(x_j | x_{j-1}) (\mu_{j \leftarrow (j+1)}(x_j) P(y_j | x_j))$$

F4



## Sum-Product Algorithm for HMMs



- Backward messages: Exactly the same, just reverse time arrow

① Measurement: Prior  $\rightarrow$  posterior:  $\propto \mu_{j \leftarrow (j+1)}(x_j) P(y_j | x_j)$

② Diffusion, information propagation (marginalization):

$$\mu_{(j-1) \leftarrow j}(x_{j-1}) \propto \sum_{x_j} P(x_j | x_{j-1}) (\mu_{j \leftarrow (j+1)}(x_j) P(y_j | x_j))$$

- Forward / backward pass independent: can be run in parallel
- Posterior marginals:

$$P(x_j | \mathbf{y}_{1 \dots J}) \propto \mu_{(j-1) \rightarrow j}(x_j) P(y_j | x_j) \mu_{j \leftarrow (j+1)}(x_j)$$

$$P(x_{j-1}, x_j | \mathbf{y}_{1 \dots J}) \propto \mu_{(j-2) \rightarrow (j-1)}(x_{j-1}) P(y_{j-1} | x_{j-1}) P(x_j | x_{j-1}) P(y_j | x_j) \mu_{j \leftarrow (j+1)}(x_j)$$

F4b

# HMMs in Practice

Enormously influential, both in practice and algorithm development

Speech recognition:

SR today: = **HMMs** with clever search tricks

- EM came from there (Baum, Welch: Forward-backward algorithm)
- Swiped field clean of anything else (rule-based, hand-coded, linguistic, ...) in 1970s. Early work at CMU (Baker, Lowerre) and IBM (Jelinek)
- $x_j$ : Subphonemes.  $\mathbf{y}_j$ : Spectral features of acoustic waveform.  
 $P(\mathbf{y}|x)$ : Gaussian mixture
- One of the big success stories of statistical learning over other “loftier” approaches
- Today more industry than research: Big groups, big computers, huge amounts of data

# HMMs in Practice

Enormously influential, both in practice and algorithm development

Bio-Informatics:

Introduced there early 90s by David Haussler (machine learning theorist, turned famous computational biologist)

- Before that: Dynamic programming sequence alignment (BLAST)
- Most macromolecules of organic chemistry **are** chains (some folded in complex ways):  
 $x_j \in \{A, C, T, G\}$  (or triplets),  $x_j \in \{\text{amino acids}\}$
- Sequence matching by pair HMMs (two  $y_j$  chains, common  $x_j$ )
- Gene finding: HGP estimates about # human genes: HMMs
- Protein categorization (homologues)
- Together with tree models: Phylogenetics, evolutionary history of species

# Further Remarks

- How do I obtain  $\operatorname{argmax}_{\{\mathbf{x}_j\}} \log P(\{\mathbf{y}_j\}, \{\mathbf{x}_j\})$ ?

For example: How to I decode words from acoustic waveform?

⇒ Max-product algorithm: **Viterbi** decoding

- Learning with inner (Viterbi) maximization usually used in SR:  
Faster than EM (beam search, pruning)

# Further Remarks

- How do I obtain  $\operatorname{argmax}_{\{\mathbf{x}_j\}} \log P(\{\mathbf{y}_j\}, \{\mathbf{x}_j\})$ ?  
For example: How to I decode words from acoustic waveform?  
⇒ Max-product algorithm: **Viterbi** decoding
  - Learning with inner (Viterbi) maximization usually used in SR:  
Faster than EM (beam search, pruning)
- Simple modifications give nonstationary transition kernels  
(e.g., for non-geometric duration time)

# Further Remarks

- How do I obtain  $\operatorname{argmax}_{\{\mathbf{x}_j\}} \log P(\{\mathbf{y}_j\}, \{\mathbf{x}_j\})$ ?  
For example: How to I decode words from acoustic waveform?  
⇒ Max-product algorithm: **Viterbi** decoding
  - Learning with inner (Viterbi) maximization usually used in SR:  
Faster than EM (beam search, pruning)
- Simple modifications give nonstationary transition kernels (e.g., for non-geometric duration time)
- Not hard to do hierarchical / multi-scale HMMs: Nodes expand into HMMs themselves (of potentially different length)

# Further Remarks

- How do I obtain  $\operatorname{argmax}_{\{\mathbf{x}_j\}} \log P(\{\mathbf{y}_j\}, \{\mathbf{x}_j\})$ ?  
For example: How to I decode words from acoustic waveform?  
⇒ Max-product algorithm: **Viterbi** decoding
  - Learning with inner (Viterbi) maximization usually used in SR:  
Faster than EM (beam search, pruning)
- Simple modifications give nonstationary transition kernels (e.g., for non-geometric duration time)
- Not hard to do hierarchical / multi-scale HMMs: Nodes expand into HMMs themselves (of potentially different length)
- HMMs for large state spaces: Model structure within  $\mathbf{x}_j$  itself
  - Factorial HMM [next lecture]

# Further Remarks

- How do I obtain  $\operatorname{argmax}_{\{\mathbf{x}_j\}} \log P(\{\mathbf{y}_j\}, \{\mathbf{x}_j\})$ ?  
For example: How to I decode words from acoustic waveform?  
⇒ Max-product algorithm: **Viterbi** decoding
  - Learning with inner (Viterbi) maximization usually used in SR:  
Faster than EM (beam search, pruning)
- Simple modifications give nonstationary transition kernels (e.g., for non-geometric duration time)
- Not hard to do hierarchical / multi-scale HMMs: Nodes expand into HMMs themselves (of potentially different length)
- HMMs for large state spaces: Model structure within  $\mathbf{x}_j$  itself
  - Factorial HMM [next lecture]
- Common theme: On some level: Markov chain (usually latent).  
Belief propagation
  - Message sizes independent of sequence length
  - Running time linear in sequence length



# Dynamical Systems

- The world is not discrete. Problems involving motion, co-occurrence: Differential equations, **continuous** variables
- Reasoning about uncertainty in such problems:  
**Dynamical state space models**

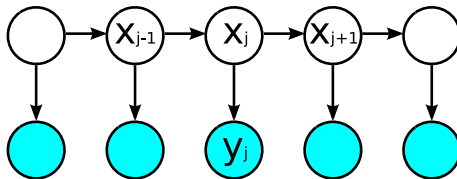
# Dynamical Systems

- The world is not discrete. Problems involving motion, co-occurrence: Differential equations, **continuous** variables
- Reasoning about uncertainty in such problems:  
**Dynamical state space models**
- **State**: Finite set variables, containing all information to move on (separate past from future). Can include derivatives as well (location, orientation, velocity, angular velocity, acceleration, torque, . . .). Usually (partly) latent
- Reason about **distribution** over state, conditioned on past observations (filtering) or all observations (smoothing). Propagate such state distributions (belief states)

# Dynamical Systems

- The world is not discrete. Problems involving motion, co-occurrence: Differential equations, **continuous** variables
  - Reasoning about uncertainty in such problems:  
**Dynamical state space models**
  - **State**: Finite set variables, containing all information to move on (separate past from future). Can include derivatives as well (location, orientation, velocity, angular velocity, acceleration, torque, ...). Usually (partly) latent
  - Reason about **distribution** over state, conditioned on past observations (filtering) or all observations (smoothing). Propagate such state distributions (belief states)
  - Use this inference for higher order tasks
    - Learning about environment, world model, sensor accuracy
    - Planning behaviour, interaction which modifies environment
- ⇒ Whatever you do: **Inference is at the bottom**

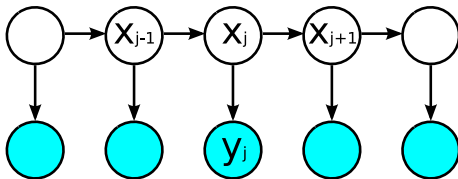
# Linear Dynamical System



- Example: Moving robot localization

F8

# Linear Dynamical System



- Example: Moving robot localization
- Local conditional probabilities: Linear-Gaussian

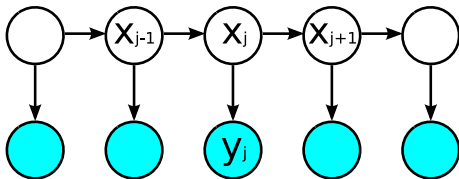
$$\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{G}_t \varepsilon_{1,t}, \quad \varepsilon_{1,t} \sim N(\mathbf{0}, \mathbf{I}) \quad \text{Transition prior}$$

$$\mathbf{y}_t = \mathbf{C}_t \mathbf{x}_t + \varepsilon_{2,t}, \quad \varepsilon_{2,t} \sim N(\mathbf{0}, \Psi_t) \quad \text{Observation likelihood}$$

$\varepsilon_{k,t}$  independent of others.

Stationary LDS:  $\mathbf{A}$ ,  $\mathbf{G}$ ,  $\mathbf{C}$ ,  $\Psi$  independent of  $t$

# Linear Dynamical System



- Example: Moving robot localization
- Local conditional probabilities: Linear-Gaussian

$$\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{G}_t \varepsilon_{1,t}, \quad \varepsilon_{1,t} \sim N(\mathbf{0}, \mathbf{I}) \quad \text{Transition prior}$$

$$\mathbf{y}_t = \mathbf{C}_t \mathbf{x}_t + \varepsilon_{2,t}, \quad \varepsilon_{2,t} \sim N(\mathbf{0}, \Psi_t) \quad \text{Observation likelihood}$$

$\varepsilon_{k,t}$  independent of others.

Stationary LDS:  $\mathbf{A}$ ,  $\mathbf{G}$ ,  $\mathbf{C}$ ,  $\Psi$  independent of  $t$

- Inference in such a model? Combine what you know:
  - HMM for discrete latent states
  - Factor analysis for linear-Gaussian model, independent states

# Behind the Equations

- Behind the equation mess in engineering textbooks, there is a **simple idea**. What you need to remember:
  - Above formulae: That idea, and the generic primitives it requires
  - Below formulae: Numerically uncritical way of implementation

# Behind the Equations

- Behind the equation mess in engineering textbooks, there is a **simple idea**. What you need to remember:
  - Above formulae: That idea, and the generic primitives it requires
  - Below formulae: Numerically uncritical way of implementation
- Simple idea:
  - All distributions / messages in this model: Gaussian.  
You need to maintain / pass:
    - Moment parameters (mean, covariance), or
    - Natural parameters
  - Belief propagation primitives: Measurement (prior  $\rightarrow$  posterior), information propagation (marginalization)



# Behind the Equations

- Behind the equation mess in engineering textbooks, there is a **simple idea**. What you need to remember:
  - Above formulae: That idea, and the generic primitives it requires
  - Below formulae: Numerically uncritical way of implementation
- Simple idea:
  - All distributions / messages in this model: Gaussian.  
You need to maintain / pass:
    - Moment parameters (mean, covariance), or
    - Natural parameters
  - Belief propagation primitives: Measurement (prior  $\rightarrow$  posterior), information propagation (marginalization)
- Below formulae: **Get these primitives right**
  - Use linear algebra to get them in the right form
  - Use numerically trusted solutions for elementary steps

# Kalman Filtering

Filtering with moment parameters:

$$\begin{array}{c}
 N(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}) \xrightarrow{\text{info. prop.}} N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \\
 \xrightarrow{\text{measurement}} N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})
 \end{array}$$

- Information propagation:  $\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{G}_t \varepsilon_{1,t}$

# Kalman Filtering

Filtering with moment parameters:

$$\begin{array}{c}
 N(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1}) \xrightarrow{\text{info. prop.}} N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \\
 \xrightarrow{\text{measurement}} N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})
 \end{array}$$

- Information propagation:  $\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{G}_t \varepsilon_{1,t}$

$$\boldsymbol{\mu}_{t|t-1} = \mathbf{A}_t \boldsymbol{\mu}_{t-1|t-1}, \quad \boldsymbol{\Sigma}_{t|t-1} = \mathbf{A}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{A}_t^T + \mathbf{G}_t \mathbf{G}_t^T$$

- Measurement: “Prior”  $N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})$ .  
Likelihood  $N(\mathbf{y}_t | \mathbf{C}_t \mathbf{x}_t, \boldsymbol{\Psi}_t)$ . Posterior?

F10

# Kalman Filtering

Filtering with moment parameters:

$$\begin{array}{c} \text{measurement} \\ \longrightarrow \end{array} N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) \xrightarrow{\text{info. prop.}} N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})$$

- Information propagation:  $\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{G}_t \boldsymbol{\varepsilon}_{1,t}$

$$\boldsymbol{\mu}_{t|t-1} = \mathbf{A}_t \boldsymbol{\mu}_{t-1|t-1}, \quad \boldsymbol{\Sigma}_{t|t-1} = \mathbf{A}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{A}_t^T + \mathbf{G}_t \mathbf{G}_t^T$$

- Measurement: “Prior”  $N(\mathbf{x}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})$ .  
Likelihood  $N(\mathbf{y}_t | \mathbf{C}_t \mathbf{x}_t, \boldsymbol{\Psi}_t)$ . Posterior?

$$\boldsymbol{\Sigma}_{t|t} = \text{Cov}[(\mathbf{x}_t \ \mathbf{y}_t)] / \text{Cov}[\mathbf{y}_t] = \boldsymbol{\Sigma}_{t|t-1} - \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}_t^T \mathbf{E}_t^{-1} \mathbf{C}_t \boldsymbol{\Sigma}_{t|t-1}$$

$$\boldsymbol{\mu}_{t|t} = \text{E}[\mathbf{x}_t] + \text{Cov}[\mathbf{x}_t, \mathbf{y}_t] \text{Cov}[\mathbf{y}_t]^{-1} (\mathbf{y}_t - \text{E}[\mathbf{y}_t])$$

$$= \boldsymbol{\mu}_{t|t-1} + \underbrace{\boldsymbol{\Sigma}_{t|t-1} \mathbf{C}_t^T \mathbf{E}_t^{-1}}_{\text{Kalman gain}} (\mathbf{y}_t - \mathbf{C}_t \boldsymbol{\mu}_{t|t-1}), \quad \mathbf{E}_t = \boldsymbol{\Psi}_t + \mathbf{C}_t \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}_t^T$$

# Remarks

- Kalman gain matrix:

$$\mathbf{K}_t = \Sigma_{t|t-1} \mathbf{C}_t^T (\Psi_t + \mathbf{C}_t \Sigma_{t|t-1} \mathbf{C}_t^T)^{-1} \left\{ = \text{Cov}[\mathbf{x}_t, \mathbf{y}_t] \text{Cov}[\mathbf{y}_t]^{-1} \right\}$$

Residual error  $\mathbf{y}_t - \mathbb{E}[\mathbf{y}_t | \mathbf{y}_{<t}] \rightarrow$  correction mean estimate.

Can also write:  $\Sigma_{t|t} = (\mathbf{I} - \mathbf{K}_t \mathbf{C}_t) \Sigma_{t|t-1}$

# Remarks

- Kalman gain matrix:

$$\mathbf{K}_t = \Sigma_{t|t-1} \mathbf{C}_t^T (\Psi_t + \mathbf{C}_t \Sigma_{t|t-1} \mathbf{C}_t^T)^{-1} \left\{ = \text{Cov}[\mathbf{x}_t, \mathbf{y}_t] \text{Cov}[\mathbf{y}_t]^{-1} \right\}$$

Residual error  $\mathbf{y}_t - \mathbb{E}[\mathbf{y}_t | \mathbf{y}_{<t}] \rightarrow$  correction mean estimate.

Can also write:  $\Sigma_{t|t} = (\mathbf{I} - \mathbf{K}_t \mathbf{C}_t) \Sigma_{t|t-1}$

- Recall: In the moment parameterization:

Information propagation Simple

Measurement Difficult (needs matrix factorization)

In the natural parameterization, these roles are reversed

- **Information filter**: Propagate natural parameters  $\mathbf{r}_{t|t}$ ,  $\mathbf{S}_{t|t}$  instead of moment parameters [ $\mathbf{S}_{t|t} = \Sigma_{t|t}^{-1}$ ,  $\mathbf{r}_{t|t} = \Sigma_{t|t}^{-1} \boldsymbol{\mu}_{t|t}$ ]

# In Practice

“Below the formulae”: What’s he talking about?

⇒ **In practice, neither of them work** (on real problems)

- In theory,  $\Sigma_{t|t}$  or  $\mathbf{S}_{t|t}$  stay positive definite. In practice they don't!
- Root of problem: Information propagation / measurement simple linear in **different** parameterizations. Conversion (matrix inversion) prone to numerical errors

# In Practice

“Below the formulae”: What’s he talking about?

⇒ **In practice, neither of them work** (on real problems)

- In theory,  $\Sigma_{t|t}$  or  $\mathbf{S}_{t|t}$  stay positive definite. In practice they don’t!
- Root of problem: Information propagation / measurement simple linear in **different** parameterizations. Conversion (matrix inversion) prone to numerical errors
- First improvement: Propagate matrix factorization:
  - Kalman square root filter  $\Sigma_{t|t} = \mathbf{F}_{t|t} \mathbf{F}_{t|t}^T$ . Propagate  $\mathbf{F}_{t|t}$
  - Information square root filter  $\Sigma_{t|t}^{-1} = \mathbf{F}_{t|t} \mathbf{F}_{t|t}^T$ . Propagate  $\mathbf{F}_{t|t}$
- Further improvements: Formulate as weighted least squares problem. Use stable LS method from numerical mathematics

[Paige, Saunders: Least Squares Estimation of Discrete Linear Dynamic Systems Using Orthogonal Transformations (1977)]



# Smoothing

- All approaches use filtering for forward pass
- Rauch-Tung-Striebel (RTS) smoother: Backward pass computes marginals  $E[\mathbf{x}_t|D]$ ,  $\text{Cov}[\mathbf{x}_t|D]$  directly,  $D = \{\mathbf{y}_t\}$ . Idea:

$$\begin{aligned}
 P(\mathbf{x}_t|D) &= \int P(\mathbf{x}_t|\mathbf{x}_{t+1}, D)P(\mathbf{x}_{t+1}|D) d\mathbf{x}_{t+1} \\
 &\stackrel{!}{=} \int P(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}_{\leq t})P(\mathbf{x}_{t+1}|D) d\mathbf{x}_{t+1} \quad [\mathbf{x}_t \perp \mathbf{y}_{>t} | \mathbf{x}_{t+1}]
 \end{aligned}$$

Work out moments of  $P(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}_{\leq t})$  from filtering variables.  
Average  $\mathbf{x}_{t+1}$  over  $P(\mathbf{x}_{t+1}|D)$ . Details: In your exercises

# Smoothing

- All approaches use filtering for forward pass
- Rauch-Tung-Striebel (RTS) smoother: Backward pass computes marginals  $E[\mathbf{x}_t|D]$ ,  $\text{Cov}[\mathbf{x}_t|D]$  directly,  $D = \{\mathbf{y}_t\}$ . Idea:

$$\begin{aligned}
 P(\mathbf{x}_t|D) &= \int P(\mathbf{x}_t|\mathbf{x}_{t+1}, D)P(\mathbf{x}_{t+1}|D) d\mathbf{x}_{t+1} \\
 &\stackrel{!}{=} \int P(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}_{\leq t})P(\mathbf{x}_{t+1}|D) d\mathbf{x}_{t+1} \quad [\mathbf{x}_t \perp \mathbf{y}_{>t} | \mathbf{x}_{t+1}]
 \end{aligned}$$

Work out moments of  $P(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}_{\leq t})$  from filtering variables.  
Average  $\mathbf{x}_{t+1}$  over  $P(\mathbf{x}_{t+1}|D)$ . Details: In your exercises

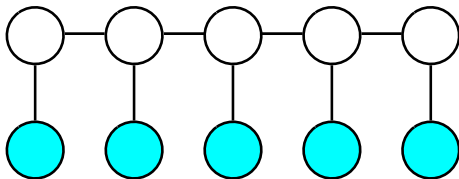
- Two-filter smoothing: Analogous to forward-backward BP
  - Run backward filter (in parallel to forward filter)
  - Combine results by Gaussian product formula. Do not count observation twice!

F13

# Learning

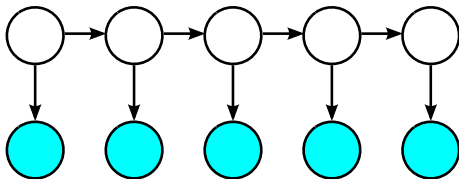
- Recall last lecture: Most difficult part of learning is inference
- EM algorithm: E step is smoothing. M step: Like in factor analysis
- Gradient-based optimization: Average gradients of log-potentials over marginal posterior (smoothing)
- Formulae even worse, but we are not impressed
  - Above formulae: Decomposition:
    - Marginal inference (smoothing)
    - Gradient accumulation, given marginals
    - Parameter updates
  - Below formulae: Use stable filtering / smoothing implementation. Gradient accumulation typically harmless

# Conditional Random Fields



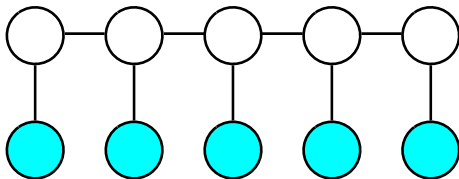
- Undirected sequence model. Different properties through global normalization (HMM: local normalization)
- Markov random field with tree graph
- Heavily used in text / language modelling (labeling, named entity recognition)
- Training with complete data (all  $\mathbf{x}_t$  given): Iterative, but convex optimization (for log-linear potentials). Can be done very efficiently (Quasi Newton optimization; approximate Newton optimization with Hessian-vector product)
- Training with incomplete data: EM outer loop required

# Conditional Random Fields



- Undirected sequence model. Different properties through global normalization (**HMM: local normalization**)
- Markov random field with tree graph
- Heavily used in text / language modelling (labeling, named entity recognition)
- Training with complete data (all  $\mathbf{x}_t$  given): Iterative, but convex optimization (for log-linear potentials). Can be done very efficiently (Quasi Newton optimization; approximate Newton optimization with Hessian-vector product)
- Training with incomplete data: EM outer loop required

# Conditional Random Fields

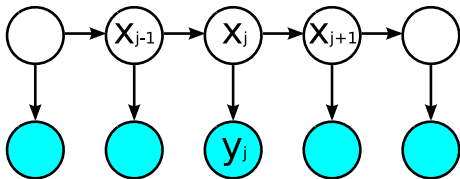


- Undirected sequence model. Different properties through global normalization (HMM: local normalization)
- Markov random field with tree graph
- Heavily used in text / language modelling (labeling, named entity recognition)
- Training with complete data (all  $\mathbf{x}_t$  given): Iterative, but convex optimization (for log-linear potentials). Can be done very efficiently (Quasi Newton optimization; approximate Newton optimization with Hessian-vector product)
- Training with incomplete data: EM outer loop required

# General State Space Models

- We've had state space models

- All discrete:
  - Multinomial family
- All linear-Gaussian:
  - Gaussian family



- What about other situations?

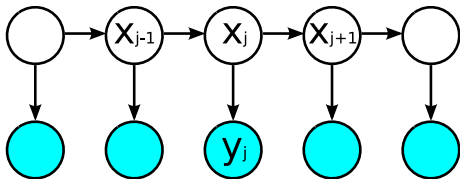
- Nonlinear dynamical system: Transition / observation mapping nonlinear. Noise dependent on current state
- Switching state space model: State consists of continuous **and** discrete variables

⇒ Graph is still a chain. Efficient inference by BP?

# General State Space Models

- We've had state space models

- All discrete:
  - Multinomial family
- All linear-Gaussian:
  - Gaussian family



- What about other situations?

- Nonlinear dynamical system: Transition / observation mapping nonlinear. Noise dependent on current state
- Switching state space model: State consists of continuous **and** discrete variables

⇒ Graph is still a chain. Efficient inference by BP?

- Problem: Only multinomial, Gaussian families
  - Closed under conditioning **and** marginalization
  - Fixed-size parameterization

⇒ Inference for general state space models: **NP hard**



# Approximate Filtering

Dynamic programming requires fixed-size message representations

⇒ Blow-up has to be countered by approximations

- Approximate transition / observation potentials locally
  - Extended Kalman filter: Linearize transition / observation mapping by Taylor expansion at  $\mu_{t-1|t-1}$  /  $\mu_{t|t-1}$  resp. Use exact propagation with approximated potentials

# Approximate Filtering

Dynamic programming requires fixed-size message representations

⇒ Blow-up has to be countered by approximations

- Approximate transition / observation potentials locally
  - Extended Kalman filter: Linearize transition / observation mapping by Taylor expansion at  $\mu_{t-1|t-1}$  /  $\mu_{t|t-1}$  resp. Use exact propagation with approximated potentials
- Use exact local propagation, project resulting message back into representation family. Projection done optimally by **matching family moments** (mean, covariance for Gaussian). Variants:
  - Assumed density filter
  - Unscented filter (cheap quadrature by exact monomials)

# Approximate Filtering

Dynamic programming requires fixed-size message representations

⇒ Blow-up has to be countered by approximations

- Approximate transition / observation potentials locally
  - Extended Kalman filter: Linearize transition / observation mapping by Taylor expansion at  $\mu_{t-1|t-1}$  /  $\mu_{t|t-1}$  resp. Use exact propagation with approximated potentials
- Use exact local propagation, project resulting message back into representation family. Projection done optimally by **matching family moments** (mean, covariance for Gaussian). Variants:
  - Assumed density filter
  - Unscented filter (cheap quadrature by exact monomials)
- EKF cheap and cheerful. ADF works better in general, but needs quadrature to approximate moments. Projection more general:
  - EKF can be seen as special backprojection as well

# What about Smoothing?

- Approximate filters: One-shot approach, local approximations never re-visited. But: Approximate inference is **iterative**
- Correct learning requires marginals given all data (also future)  
⇒ Smoothing

# What about Smoothing?

- Approximate filters: One-shot approach, local approximations never re-visited. But: Approximate inference is **iterative**
- Correct learning requires marginals given all data (also future)  
⇒ Smoothing
- Options:
  - Apply filter approximation technique to smoother (e.g., two-filter smoother) ⇒ not iterative
  - Better: Use principled approximate inference framework [expectation propagation, part II]
- General warning: Numerically even more difficult than inference in fixed LDS ⇒ Attention to numerical details essential

# Wrap-Up

- Hidden Markov model (discrete states): Non-Markovian behaviour from Markovian ingredients
- Linear dynamical system: Simple idea, messy equations. Does not work without numerically careful implementation
- Conditional random field: Alternative to HMM for large text / language problems
- Filtering / smoothing for general state space models: Approximation by (moment matching) backprojection