

# EM Algorithm for Learning Linear Dynamical Systems

May.10, 2018

Takehisa YAIRI (矢入健久)

E-mail: [yairi@ailab.t.u-tokyo.ac.jp](mailto:yairi@ailab.t.u-tokyo.ac.jp)

# Today's Goal

- Understand the maximum likelihood approach to learning dynamical systems
- Derive the EM (Expectation-Maximization) algorithm for linear dynamical systems
- Need to overcome many hurdles to reach there
  - Bayesian inference for dynamical systems
  - Properties of linear Gaussian systems
  - Kalman filtering and RTS smoothing
  - Jensen's inequality
  - Minorization-maximization (MM) algorithm
- It's a long journey...



# Probabilistic Inference for Dynamical Systems

- *Introduction of Bayesian Filtering* -

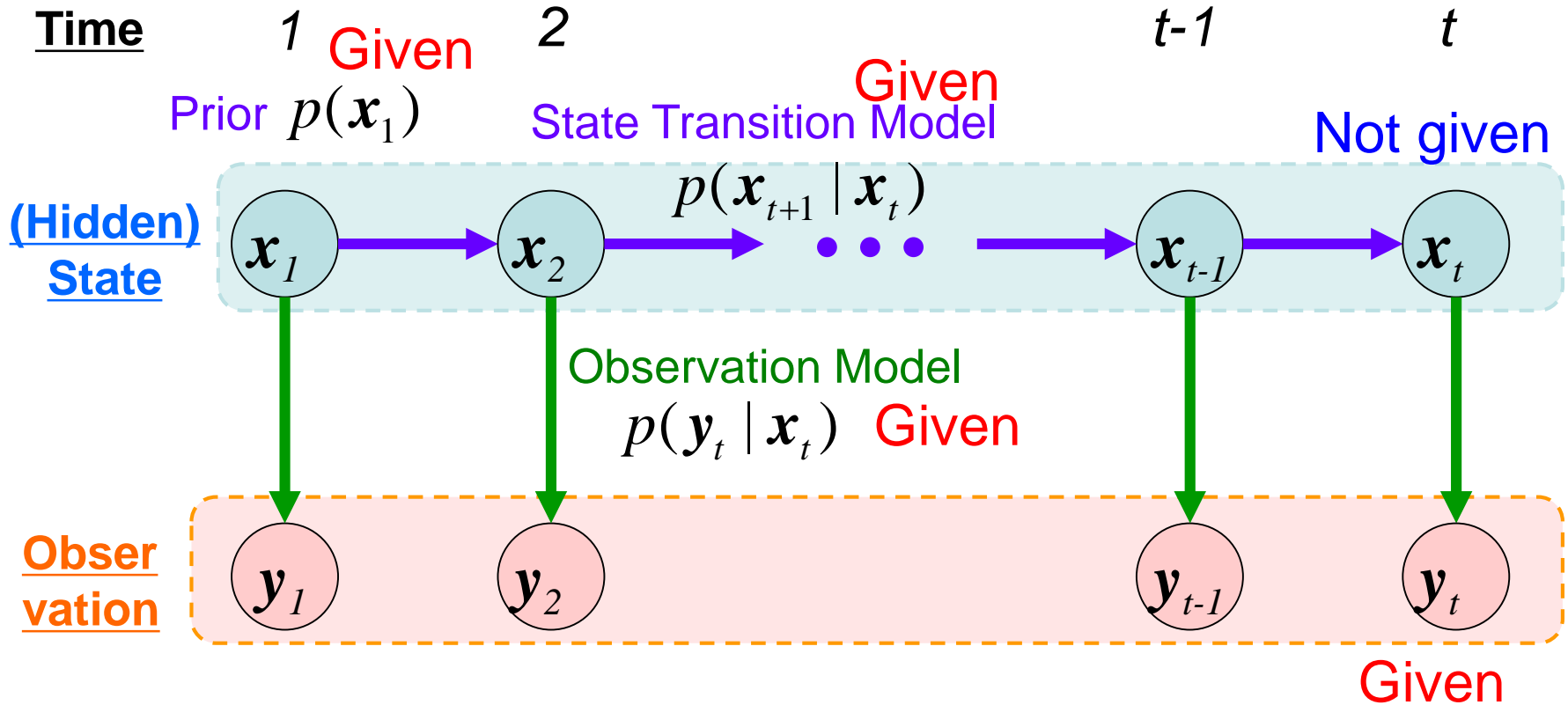
# Bayesian Filtering

## (Problem Definition)

- Purpose :
  - Estimate the current state  $\mathbf{x}_t$  optimally using all **observations up to now**  $\mathbf{y}_{1:t}$
- Given :
  - State transition model :  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$
  - Observation model :  $p(\mathbf{y}_t | \mathbf{x}_t)$
  - Prior distribution :  $p(\mathbf{x}_1)$
  - Observations up to  $t$  :  $\mathbf{y}_{1:t} \equiv \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{t-1}, \mathbf{y}_t\}$
- Find :
  - Posterior distribution at time  $t$  :  $p(\mathbf{x}_t | \mathbf{y}_{1:t})$

For simplicity,  
control input  $\mathbf{u}_t$   
is ignored

# Illustration



## State Estimation:

When    and → and ↓ are given, estimate   

observation
State model
Observation model
Hidden state

# Bayesian Filtering: Derivation (1)

All observation **up to time t**      Obs. at time **t**      Obs. **up to t-1**

$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) = p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{y}_{t-1}, \dots, \mathbf{y}_2, \mathbf{y}_1) = p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{y}_{1:t-1})$$

$$= \frac{p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{y}_{1:t-1}) \cdot p(\mathbf{x}_t | \mathbf{y}_{1:t-1})}{p(\mathbf{y}_t | \mathbf{y}_{1:t-1})}$$

Bayes' theorem  $P(X | Y) = \frac{P(Y | X) \cdot P(X)}{P(Y)}$

Independent of  $\mathbf{x}_t$

$$= \alpha \cdot p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{y}_{1:t-1}) \cdot p(\mathbf{x}_t | \mathbf{y}_{1:t-1})$$

constant      Predictive distribution

# Bayesian Filtering: Derivation (2)

$$p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{y}_{1:t-1}) = p(\mathbf{y}_t | \mathbf{x}_t) \quad \leftarrow \begin{array}{l} \text{Markov property} \\ \text{(Conditional independence)} \end{array}$$

Observation model

$$p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{x}_t, \mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} \quad \leftarrow \text{Marginalization}$$

$$= \int [p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}_{1:t-1}) \cdot p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1})] d\mathbf{x}_{t-1} \quad \leftarrow \begin{array}{l} \text{Bayes' Rule} \end{array}$$

$$\approx \int \underbrace{p(\mathbf{x}_t | \mathbf{x}_{t-1})}_{\text{State Transition model}} \cdot \underbrace{p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1})}_{\text{Posterior at previous time step}} d\mathbf{x}_{t-1} \quad \leftarrow \begin{array}{l} \text{Markov property} \end{array}$$

# Bayesian Filtering: Derivation (3)

After all,

$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) \approx \alpha \cdot p(\mathbf{y}_t | \mathbf{x}_t) \int [p(\mathbf{x}_t | \mathbf{x}_{t-1}) \cdot p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1})] d\mathbf{x}_{t-1}$$

Posterior (Belief) at time t      Observation Model      Transition Model      Posterior (Belief) at time t-1

- Posterior at every time step can be computed **recursively**
  - On-line update
  - No need to store past measurements
- What should we do with  $\alpha$  (normalizing factor) ?
  - Automatically determined so that posterior integrated over  $X_t$  become 1



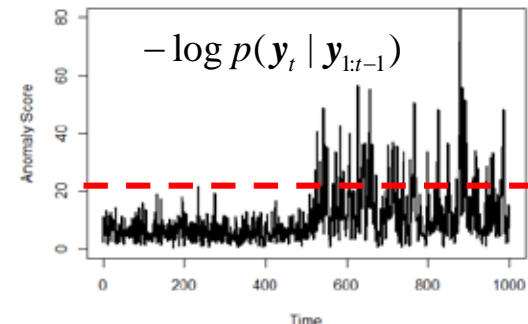
# Normalizing Factor

The normalizing factor  $\alpha$  can be computed explicitly, though we ignored it as it is independent of  $\mathbf{x}_t$

$$\begin{aligned}\alpha^{-1} &= p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{y}_t, \mathbf{x}_t | \mathbf{y}_{1:t-1}) d\mathbf{x}_t \\ &= \int p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{y}_{1:t-1}) \cdot p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) d\mathbf{x}_t \\ &= \int p(\mathbf{y}_t | \mathbf{x}_t) \cdot \left[ \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} \right] d\mathbf{x}_t\end{aligned}$$

Predictive distribution for the next observation

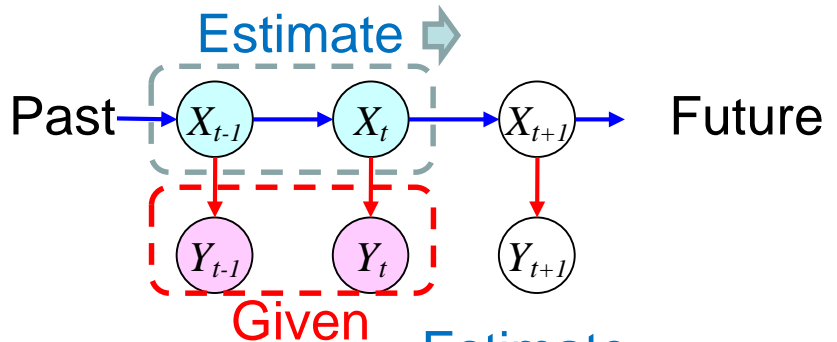
➡ Useful for anomaly detection !



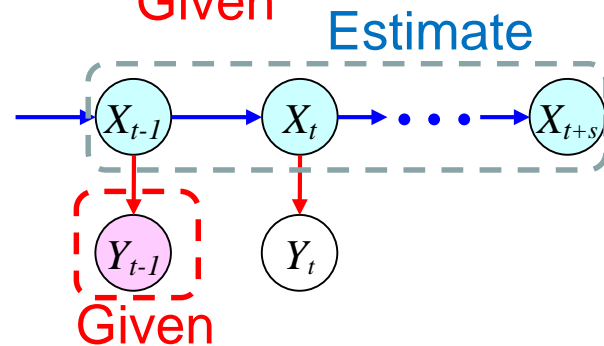
# Prediction and Smoothing

- “Prediction” and “smoothing” are similar to but different from “filtering”

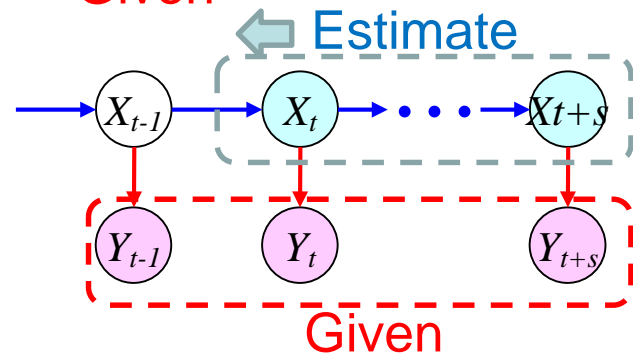
Filtering :  $p(\mathbf{x}_t | \mathbf{y}_{1:t})$



Prediction :  $p(\mathbf{x}_{t+s} | \mathbf{y}_{1:t})$

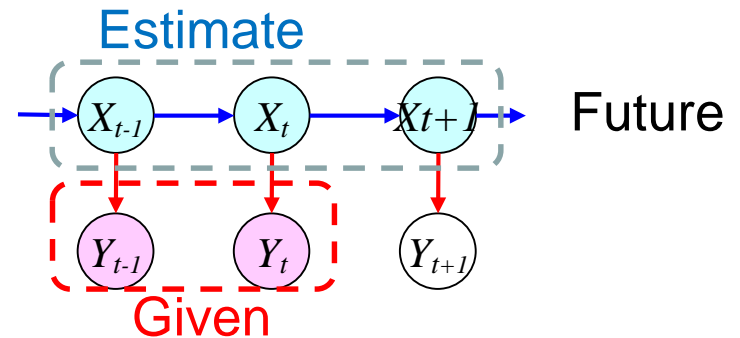


Smoothing :  $p(\mathbf{x}_t | \mathbf{y}_{1:t+s})$



# Bayesian Prediction (1)

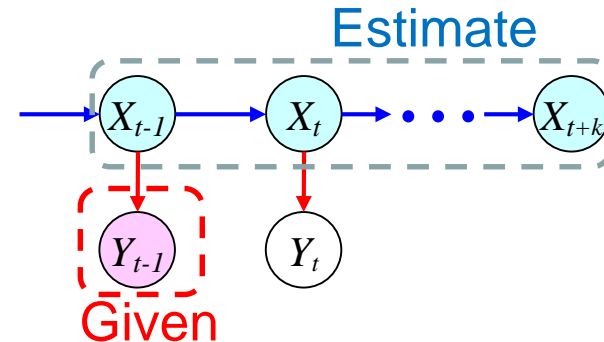
One-step Prediction



$$\begin{aligned}
 p(\mathbf{x}_{t+1} / \mathbf{y}_{1:t}) &= \int p(\mathbf{x}_{t+1}, \mathbf{x}_t / \mathbf{y}_{1:t}) d\mathbf{x}_t && \text{Marginalization} \\
 &= \int p(\mathbf{x}_{t+1} | \mathbf{x}_t, \cancel{\mathbf{y}_{1:t}}) p(\mathbf{x}_t | \mathbf{y}_{1:t}) d\mathbf{x}_t && \text{Bayes rule} \\
 &= \int \underbrace{p(\mathbf{x}_{t+1} | \mathbf{x}_t)}_{\substack{\text{State} \\ \text{Transition} \\ \text{model}}} \underbrace{p(\mathbf{x}_t | \mathbf{y}_{1:t})}_{\substack{\text{Filtering} \\ \text{Distribution} \\ \text{at time } t}} d\mathbf{x}_t && \text{Markov property}
 \end{aligned}$$

# Bayesian Prediction (2)

***k*-step Prediction**



$$\begin{aligned}
 p(\mathbf{x}_{t+k} \mid \mathbf{y}_{1:t}) &= \int p(\mathbf{x}_{t+k}, \mathbf{x}_{t+k-1} \mid \mathbf{y}_{1:t}) d\mathbf{x}_{t+k-1} \\
 &= \iint p(\mathbf{x}_{t+k}, \mathbf{x}_{t+k-1}, \mathbf{x}_{t+k-2} \mid \mathbf{y}_{1:t}) d\mathbf{x}_{t+k-1} d\mathbf{x}_{t+k-2} \\
 &= \iiint \cdots \int p(\mathbf{x}_{t+k}, \mathbf{x}_{t+k-1}, \cdots, \mathbf{x}_t \mid \mathbf{y}_{1:t}) d\mathbf{x}_{t+k-1} \cdots d\mathbf{x}_t
 \end{aligned}$$

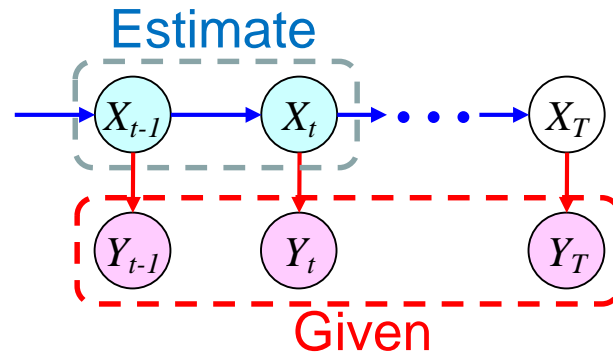
**Multiple integral**

$$= \underbrace{\int p(\mathbf{x}_{t+k} \mid \mathbf{x}_{t+k-1})}_{\text{State Transition model}} \underbrace{p(\mathbf{x}_{t+k-1} \mid \mathbf{y}_{1:t})}_{\text{(k-1)-step prediction}} d\mathbf{x}_{t+k-1} \quad \text{Recursive form}$$

# Bayesian Smoothing (1)

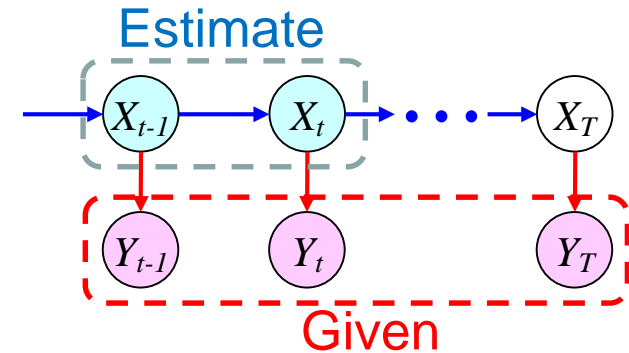
- Derive the equation for computing the posterior of state variable in smoothing where  $t < T$

$T$  : terminal time step



- Consider a backward recursive form !
  - i.e., Represent  $p(\mathbf{x}_t | \mathbf{y}_{1:T})$  in terms of  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:T})$

# Bayesian Smoothing (2)



$$p(\mathbf{x}_t \mid \mathbf{y}_{1:T}) = \int p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) d\mathbf{x}_{t+1}$$

marginal of joint distribution with  $\mathbf{x}_{t+1}$

$$p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) = p(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{y}_{1:T}) \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) \quad \text{Bayes rule}$$

$$= p(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{y}_{1:t}) \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) \quad \text{Markov property}$$

$$= \frac{p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:t})}{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t})} \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) \quad \leftarrow P(A|B) = P(A,B)/P(B)$$

$$= \frac{p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \mathbf{y}_{1:t}) \cdot p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T})}{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t})}$$

# Bayesian Smoothing (3)

$$\begin{aligned}
 p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) &= \frac{p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \cancel{\mathbf{y}_{1:t}}) \cdot p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T})}{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t})} \\
 \text{Joint} \\
 &= \frac{p(\mathbf{x}_{t+1} \mid \mathbf{x}_t) \cdot p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T})}{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t})} \quad \text{Markov property}
 \end{aligned}$$

Finally,


$$\begin{aligned}
 &\text{Smoothing at time } t \\
 p(\mathbf{x}_t \mid \mathbf{y}_{1:T}) &= \int \underbrace{p(\mathbf{x}_{t+1} \mid \mathbf{x}_t)}_{\text{Transition model}} \cdot \underbrace{p(\mathbf{x}_t \mid \mathbf{y}_{1:t})}_{\text{Filtering at time } t} \cdot \underbrace{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T})}_{\text{Smoothing at } t+1} \underbrace{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t})}_{\text{One-step prediction at } t} d\mathbf{x}_{t+1}
 \end{aligned}$$

# Properties of Linear Gaussian Systems



# Multivariate Gaussian Distribution (1)

Consider that vectors  $\mathbf{x}$ ,  $\mathbf{y}$  are multivariate Gaussian

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) \quad \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$$


Covariance matrices

$$\boldsymbol{\Sigma}_x = V(\mathbf{x}) \equiv E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T]$$

Cross-Covariance matrix between  $\mathbf{x}$  and  $\mathbf{y}$

$$\boldsymbol{\Sigma}_{xy} = \text{Cov}(\mathbf{x}, \mathbf{y}) \equiv E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^T]$$

Joint distribution is also a multivariate Gaussian

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{xy}^T & \boldsymbol{\Sigma}_y \end{bmatrix}\right) \quad (\text{Joint Gaussian})$$

# Multivariate Gaussian Distribution (2)

- Linear combination

$x$  and  $y$  are independent

Assume  $\mathbf{x} \sim N(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$   $\mathbf{y} \sim N(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$   $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbf{0}$

Then  $A\mathbf{x} + B\mathbf{y} \sim N(A\boldsymbol{\mu}_x + B\boldsymbol{\mu}_y, A\boldsymbol{\Sigma}_x A^T + B\boldsymbol{\Sigma}_y B^T)$

- If  $\mathbf{x}$  and  $\mathbf{y}$  are joint Gaussian (and mutually dependent), then conditional probability distribution of  $\mathbf{y}$  given  $\mathbf{x}$  (or posterior of  $\mathbf{y}$ ) is obtained as:

$$p(\mathbf{y} | \mathbf{x}) = N\left(\mathbf{y} | \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{xy}^T \boldsymbol{\Sigma}_x^{-1} (\mathbf{x} - \boldsymbol{\mu}_x), \boldsymbol{\Sigma}_y - \boldsymbol{\Sigma}_{xy}^T \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\Sigma}_{xy}\right)$$

cf. Prior :  $P(\mathbf{y}) = N(\mathbf{y} | \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$

Essence of  
Kalman filter

# Important Results for Linear Gaussian Systems (1)

Assume:

$x$  is a Gaussian :  $x \sim N(\mathbf{m}, \mathbf{P}) \iff p(x) = N(x | \mathbf{m}, \mathbf{P})$

$y$  is a linear transform of  $x$  plus Gaussian noise  $v$ :

$v \sim N(\mathbf{0}, \mathbf{R}) \iff p(v) = N(v | \mathbf{0}, \mathbf{R})$

$y = \mathbf{A}x + \mathbf{b} + v \iff p(y | x) = N(y | \mathbf{A}x + \mathbf{b}, \mathbf{R})$

Joint distribution:

$$p(x, y) = p(y | x)p(x) = N(y | \mathbf{A}x + \mathbf{b}, \mathbf{R}) \cdot N(x | \mu, \mathbf{P})$$

$$= N\left(\begin{bmatrix} x \\ y \end{bmatrix} \middle| \begin{bmatrix} \mu \\ \mathbf{A}\mu + \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{P} & \mathbf{P}\mathbf{A}^T \\ \mathbf{A}\mathbf{P} & \mathbf{A}\mathbf{P}\mathbf{A}^T + \mathbf{R} \end{bmatrix}\right)$$

Marginal distribution:

$$p(y) = \int p(x, y)dx = N(y | \mathbf{A}\mu + \mathbf{b}, \mathbf{A}\mathbf{P}\mathbf{A}^T + \mathbf{R})$$

# Important Results for Linear Gaussian Systems (2)

Assume joint distribution of  $\mathbf{x}$  and  $\mathbf{y}$  is Gaussian:

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{P} & \mathbf{R} \\ \mathbf{R}^T & \mathbf{Q} \end{bmatrix}\right)$$

Marginal :  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{P}) \quad p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{b}, \mathbf{Q})$

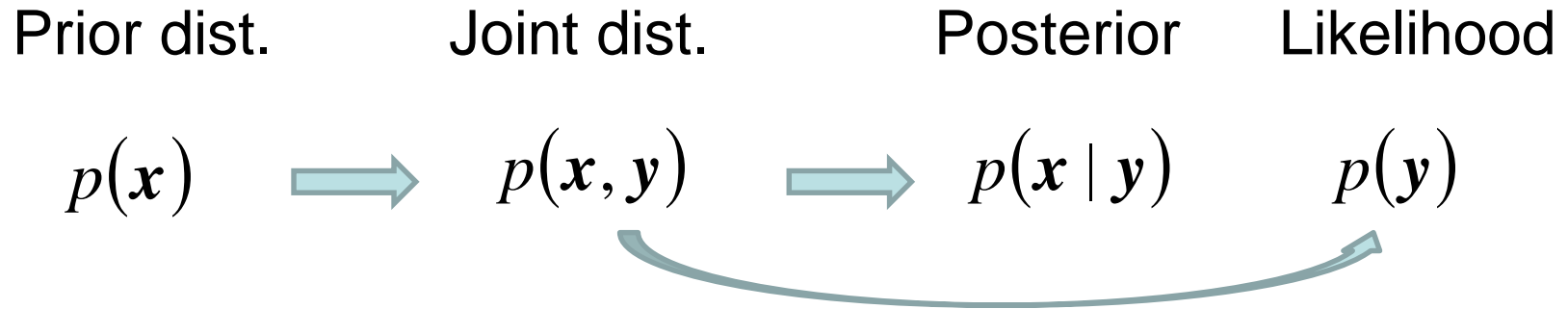
Conditional distribution:

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}\left(\mathbf{x} | \mathbf{a} + \mathbf{R}\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{b}), \mathbf{P} - \mathbf{R}\mathbf{Q}^{-1}\mathbf{R}^T\right)$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}\left(\mathbf{y} | \mathbf{b} + \mathbf{R}^T\mathbf{P}^{-1}(\mathbf{x} - \mathbf{a}), \mathbf{Q} - \mathbf{R}^T\mathbf{P}^{-1}\mathbf{R}\right)$$

# Important Results for Linear Gaussian Systems (3)

With these two results, we can compute



for linear Gaussian systems.

Imagine that  $\mathbf{x}$  is (hidden) state vector,  
and  $\mathbf{y}$  is measurement vector

# Inference for Linear Dynamical Systems

- Kalman filtering and RTS smoothing -

# Assumptions of Linear Gaussian

- Linear dynamical system  
(state space representation)

$$\begin{cases} \mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{w}_t & \text{(State equation)} \\ \mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_t & \text{(Observation equation)} \end{cases}$$

- Noises are Gaussian (and independent)

$$\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$$

- Prior distribution of initial state is Gaussian

$$p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 \mid \mathbf{m}_0, \mathbf{V}_0)$$

- Assume filter distribution  $p(\mathbf{x}_t \mid \mathbf{y}_{1:t})$  is also Gaussian

$$p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t \mid \mathbf{m}_t, \mathbf{V}_t)$$

# Derivation of Kalman Filter (1)

- Goal : Derive a recursive equation of filter dist.

$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t) \quad (\text{Filter distribution at } t)$$

$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{m}_{t+1}, \mathbf{V}_{t+1}) \quad (\text{Filter distribution at } t+1)$$



As filter distribution is Gaussian,

Derive  $\mathbf{m}_{t+1}, \mathbf{V}_{t+1}$  from  $\mathbf{m}_t, \mathbf{V}_t$

- Two steps :

1. Compute predictive distribution :  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t})$

2. Update filter distribution :  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t}, \mathbf{y}_{t+1})$



# Derivation of Kalman Filter (2): Prediction

Start with filter distribution at time  $t$

$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t)$$

State Equation

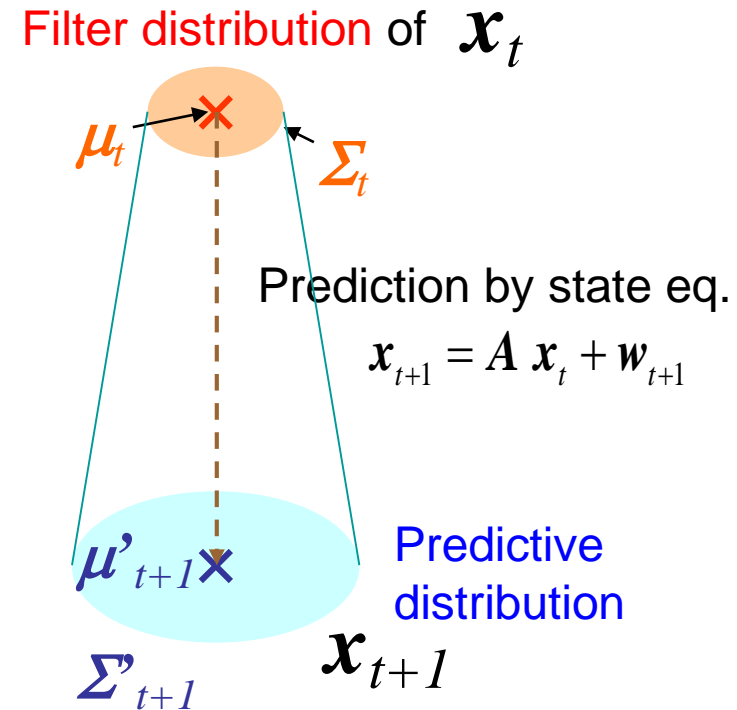
$$\begin{cases} \mathbf{x}_{t+1} = \mathbf{A} \mathbf{x}_t + \mathbf{w}_{t+1} \\ \mathbf{w}_{t+1} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \end{cases}$$

⇓ Equivalent

$$p(\mathbf{x}_{t+1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{x}_t, \mathbf{Q})$$

One-step Prediction

$$\begin{aligned} p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t}) &= \int p(\mathbf{x}_{t+1}, \mathbf{x}_t | \mathbf{y}_{1:t}) d\mathbf{x}_t \\ &= \int p(\mathbf{x}_{t+1} | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{y}_{1:t}) d\mathbf{x}_t \end{aligned}$$



How can we compute this ?

# Derivation of Kalman Filter (2): Prediction

Remember the “important property of Gaussian Distribution” !

–  $\mathbf{x}_{t+1}$  is a linear transform of  $\mathbf{x}_t$  plus Gaussian noise

Joint distribution of  $\mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  :

$$\begin{aligned} p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:t}) &= p(\mathbf{x}_{t+1} \mid \mathbf{x}_t) p(\mathbf{x}_t \mid \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_{t+1} \mid \mathbf{A}\mathbf{x}_t, \mathbf{Q}) \cdot \mathcal{N}(\mathbf{x}_t \mid \boldsymbol{\mu}_t, \mathbf{V}_t) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_t \\ \mathbf{A}\mathbf{m}_t \end{bmatrix}, \begin{bmatrix} \mathbf{V}_t & \mathbf{V}_t\mathbf{A}^T \\ \mathbf{A}\mathbf{V}_t & \mathbf{A}\mathbf{V}_t\mathbf{A}^T + \mathbf{Q} \end{bmatrix}\right) \end{aligned}$$

Marginal distribution of  $\mathbf{x}_{t+1}$  : **Predictive distribution !**

$$p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:t}) = \int p(\mathbf{x}_{t+1}, \mathbf{x}_t \mid \mathbf{y}_{1:t}) d\mathbf{x}_t = \mathcal{N}\left(\mathbf{x}_{t+1} \mid \mathbf{A}\mathbf{m}_t, \underbrace{\mathbf{A}\mathbf{V}_t\mathbf{A}^T + \mathbf{Q}}_{\substack{\text{def.} \\ \mathbf{P}_t}}\right)$$

# Derivation of Kalman Filter: Update

- Predictive dist.:  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{m}_t, \mathbf{P}_t)$
- Now that a new measurement  $\mathbf{y}_{t+1}$  comes in, we want the filter dist. at  $t+1$ :  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1})$
- Consider the joint dist. of  $\mathbf{x}_{t+1}$  and  $\mathbf{y}_{t+1}$ , given  $\mathbf{y}_{1:t}$
- Observation equation:

$$\left\{ \begin{array}{l} \mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_t \\ \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}) \end{array} \right. \quad \begin{array}{c} \text{equival.} \\ \text{equiv.} \end{array} \quad p(\mathbf{y}_{t+1} | \mathbf{x}_{t+1}) = \mathcal{N}(\mathbf{y}_{t+1} | \mathbf{C}\mathbf{x}_{t+1}, \mathbf{R})$$

$$p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} | \mathbf{y}_{1:t}) = \mathcal{N}\left(\left[\begin{array}{c} \mathbf{x}_{t+1} \\ \mathbf{y}_{t+1} \end{array}\right] \middle| \left[\begin{array}{c} \mathbf{A}\mathbf{m}_t \\ \mathbf{C}\mathbf{A}\mathbf{m}_t \end{array}\right], \left[\begin{array}{cc} \mathbf{P}_t & \mathbf{P}_t\mathbf{C}^T \\ \mathbf{C}\mathbf{P}_t & \mathbf{C}\mathbf{P}_t\mathbf{C}^T + \mathbf{R} \end{array}\right]\right)$$

$$\text{where } \mathbf{P}_t = \mathbf{A}\mathbf{V}_t\mathbf{A}^T + \mathbf{Q}$$

# Derivation of Kalman Filter: Update

- Recall another important property of Gaussian !

$$p(x, y) = N\left(\begin{bmatrix} x \\ y \end{bmatrix} \middle| \begin{bmatrix} a \\ b \end{bmatrix}, \begin{bmatrix} P & R \\ R^T & Q \end{bmatrix}\right) \quad \text{Joint}$$



$$p(x | y) = N\left(x \middle| a + RQ^{-1}(y - b), P - RQ^{-1}R^T\right) \quad \text{Conditional}$$

$$p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} | \mathbf{y}_{1:t}) = N\left(\begin{bmatrix} \mathbf{x}_{t+1} \\ \mathbf{y}_{t+1} \end{bmatrix} \middle| \begin{bmatrix} A\mathbf{m}_t \\ C A\mathbf{m}_t \end{bmatrix}, \begin{bmatrix} P_t & P_t C^T \\ C P_t & C P_t C^T + R \end{bmatrix}\right)$$



Confirm yourself !

$$\begin{aligned} p(\mathbf{x}_{t+1} | \mathbf{y}_{t+1}, \mathbf{y}_{1:t}) &= p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) \quad \text{Filter dist. at } t+1 \\ &= N(\mathbf{x}_{t+1} | A\mathbf{m}_t + K_{t+1}(\mathbf{y}_{t+1} - C A\mathbf{m}_t), (I - K_{t+1}C)P_t) \end{aligned}$$

$$\text{where } K_{t+1} = P_t C^T (C P_t C^T + R)^{-1}$$

# Summary: Linear Kalman Filter

- Input :  $p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t)$  and  $\mathbf{y}_{t+1}$
- Prediction:  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{m}_t, \mathbf{P}_t)$ 
  - where  $\mathbf{P}_t = \mathbf{A}\mathbf{V}_t\mathbf{A}^T + \mathbf{Q}$
- Update :
$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{m}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{C}\mathbf{A}\mathbf{m}_t), (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{C})\mathbf{P}_t)$$
  - where  $\mathbf{K}_{t+1} = \mathbf{P}_t\mathbf{C}^T(\mathbf{C}\mathbf{P}_t\mathbf{C}^T + \mathbf{R})^{-1}$
- Output :  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{m}_{t+1}, \mathbf{V}_{t+1})$ 
  - where 
$$\begin{cases} \mathbf{m}_{t+1} = \mathbf{A}\mathbf{m}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{C}\mathbf{A}\mathbf{m}_t) \\ \mathbf{V}_{t+1} = (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{C})\mathbf{P}_t \end{cases}$$

# Derivation of RTS Smoother (1)

- Goal : Backward recursive equation of smoothed dist.

$$p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_{t+1} \mid \hat{\mathbf{m}}_{t+1}, \hat{\mathbf{V}}_{t+1}) \quad (\text{Smoothed dist. at } t+1)$$

$$p(\mathbf{x}_t \mid \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t \mid \hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t) \quad (\text{Smoothed dist. at } t)$$



Derive  $\hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t$  from  $\hat{\mathbf{m}}_{t+1}, \hat{\mathbf{V}}_{t+1}$

# Derivation of Rauch-Tung-Striebel (RTS) Smoother \*

\* RTS smoother is also called Kalman smoother

If the joint distribution of  $\mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  given  $\mathbf{y}_{1:T}$  is obtained, we can get marginal distribution of  $\mathbf{x}_t$  using the property of Gaussian

If we know

$$p(\mathbf{x}_t, \mathbf{x}_{t+1} | \mathbf{y}_{1:T}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{bmatrix} \middle| \begin{bmatrix} \hat{\mathbf{m}}_t \\ \hat{\mathbf{m}}_{t+1} \end{bmatrix}, \begin{bmatrix} \hat{\mathbf{V}}_t & \hat{\mathbf{V}}_{t,t+1} \\ \hat{\mathbf{V}}_{t+1,t} & \hat{\mathbf{V}}_{t+1} \end{bmatrix} \right)$$

Joint Gaussian



$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t | \hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t)$$

Marginal

So, try to obtain the joint distribution  $p(\mathbf{x}_t, \mathbf{x}_{t+1} | \mathbf{y}_{1:T})$

# Derivation of RTS Smoother

$$\begin{aligned}
 p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) &= p(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{y}_{1:T}) \cdot p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T}) && \text{Bayes rule} \\
 &= \underbrace{p(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{y}_{1:t})}_{\text{Smoothed dist. at } t+1} \cdot \underbrace{p(\mathbf{x}_{t+1} \mid \mathbf{y}_{1:T})}_{\text{Markov property}}
 \end{aligned}$$

By the way,

$$p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:t}) = \mathcal{N} \left( \begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_t \\ A\mathbf{m}_t \end{bmatrix}, \begin{bmatrix} \mathbf{V}_t & \mathbf{V}_t A^T \\ A\mathbf{V}_t & \mathbf{P}_t \end{bmatrix} \right)$$

So,

$$p(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t \mid \mathbf{m}_t + \mathbf{J}_t(\mathbf{x}_{t+1} - A\mathbf{m}_t), \mathbf{V}_t - \mathbf{J}_t A \mathbf{V}_t)$$

$$\text{where } \mathbf{J}_t = \mathbf{V}_t A^T \mathbf{P}_t^{-1}$$



# Derivation of RTS Smoother

Now, we have

$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_{t+1} | \hat{\mathbf{m}}_{t+1}, \hat{\mathbf{V}}_{t+1})$$

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}_{1:T}) &= p(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}_{1:t}) \\ &= \mathcal{N}(\mathbf{x}_t | \underbrace{\mathbf{m}_t + \mathbf{J}_t(\mathbf{x}_{t+1} - \mathbf{A}\mathbf{m}_t)}_{\text{a linear transform of } \mathbf{x}_{t+1}}, \underbrace{\mathbf{V}_t - \mathbf{J}_t\mathbf{A}\mathbf{V}_t}_{\text{plus Gaussian noise}}) \end{aligned}$$

$\mathbf{x}_t$  is a linear transform of  $\mathbf{x}_{t+1}$  plus Gaussian noise

We can obtain the joint distribution !

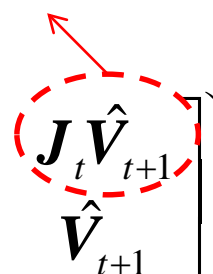
$$\begin{aligned} &p(\mathbf{x}_t, \mathbf{x}_{t+1} | \mathbf{y}_{1:T}) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_t + \mathbf{J}_t(\hat{\mathbf{m}}_{t+1} - \mathbf{A}\mathbf{m}_t) \\ \hat{\mathbf{m}}_{t+1} \end{bmatrix}, \begin{bmatrix} \mathbf{J}_t\hat{\mathbf{V}}_{t+1}\mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t\mathbf{A}\mathbf{V}_t & \mathbf{J}_t\hat{\mathbf{V}}_{t+1} \\ \hat{\mathbf{V}}_{t+1}\mathbf{J}_t^T & \hat{\mathbf{V}}_{t+1} \end{bmatrix}\right) \end{aligned}$$

# Derivation of RTS Smoother

Joint distribution

$$p(\mathbf{x}_t, \mathbf{x}_{t+1} \mid \mathbf{y}_{1:T})$$

$$= \mathcal{N} \left( \begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m}_t + \mathbf{J}_t (\hat{\mathbf{m}}_{t+1} - \mathbf{A} \mathbf{m}_t) \\ \hat{\mathbf{m}}_{t+1} \end{bmatrix}, \begin{bmatrix} \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t \mathbf{A} \mathbf{V}_t & \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \\ \hat{\mathbf{V}}_{t+1} \mathbf{J}_t^T & \hat{\mathbf{V}}_{t+1} \end{bmatrix} \right)$$

$\text{cov}[\mathbf{x}_t, \mathbf{x}_{t+1}]$   




Marginal distribution = Smoothed dist. at  $t$

$$p(\mathbf{x}_t \mid \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t \mid \hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t)$$

$$= \mathcal{N}(\mathbf{x}_t \mid \mathbf{m}_t + \mathbf{J}_t (\hat{\mathbf{m}}_{t+1} - \mathbf{A} \mathbf{m}_t), \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t \mathbf{A} \mathbf{V}_t)$$

$$\begin{cases} \hat{\mathbf{m}}_t = \mathbf{m}_t + \mathbf{J}_t (\hat{\mathbf{m}}_{t+1} - \mathbf{A} \mathbf{m}_t) \\ \hat{\mathbf{V}}_t = \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t \mathbf{A} \mathbf{V}_t = \mathbf{V}_t + \mathbf{J}_t (\hat{\mathbf{V}}_{t+1} - \mathbf{P}_t) \mathbf{J}_t^T \end{cases}$$

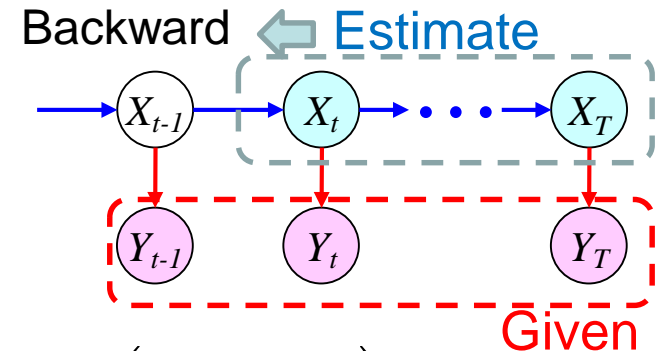
# (Summary) RTS Smoother

- Assume filtered dist.  $p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t)$  have been computed by Kalman filtering
- At terminal time,  $p(\mathbf{x}_T | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_T | \mathbf{m}_T, \mathbf{V}_T) = \mathcal{N}(\mathbf{x}_T | \hat{\mathbf{m}}_T, \hat{\mathbf{V}}_T)$
- For  $t=T-1$  to 1, repeat the following computation

$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t | \hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t)$$

where

$$\begin{cases} \hat{\mathbf{m}}_t = \mathbf{m}_t + \mathbf{J}_t (\hat{\mathbf{m}}_{t+1} - \mathbf{A} \mathbf{m}_t) \\ \hat{\mathbf{V}}_t = \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t \mathbf{A} \mathbf{V} = \mathbf{V}_t + \mathbf{J}_t (\hat{\mathbf{V}}_{t+1} - \mathbf{P}_t) \mathbf{J}_t^T \end{cases}$$

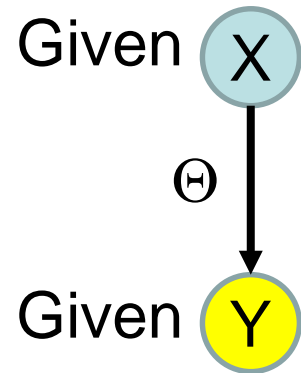


# Maximum likelihood estimation

# Case of Supervised Learning

- Assume generative models
  - $\Theta$  : parameters
- Classification (e.g., Bayesian classifier):
  - $x$ : class (categorical scalar)
  - $y$ : measurements (real / categorical vector)
- Regression (e.g., Linear Regression):
  - $x$ : input (real vector)
  - $y$ : output (real scalar/vector)
- Data :  $D = \{x_i, y_i\}$  ( $i=1, \dots, N$ ), Assume i.i.d.
- Log-likelihood function:

$$l(\Theta | D) \equiv \ln p(Y | X, \Theta) = \sum_{i=1}^N \ln p(y_i | x_i, \Theta) \Rightarrow \text{Easy to maximize}$$



# Case of Unsupervised Learning

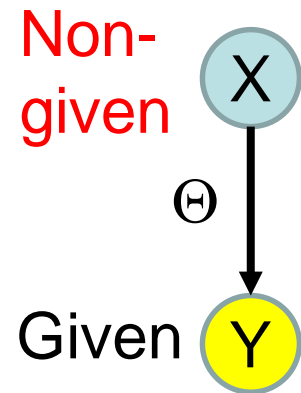
- Clustering (e.g., Gaussian mixture):
  - x: cluster (categorical scalar)
  - y: measurements (real / categorical vector)
- Dimensionality reduction
  - e.g., Factor analysis, Probabilistic PCA
  - x: factor (real vector), low-dimensional
  - y: measurements (real vector), high-dim.
- Data :  $D = \{y_i\}$  ( $i=1,\dots,N$ ), i.i.d. assumption
- Log-likelihood function:

$$\begin{aligned} l(\Theta | D) &\equiv \ln p(Y | \Theta) = \ln \int p(Y, X | \Theta) dX \\ &= \sum_{i=1}^N \ln \int p(y_i | x_i, \Theta) p(x_i | \Theta) dx_i \end{aligned}$$

Complicated !



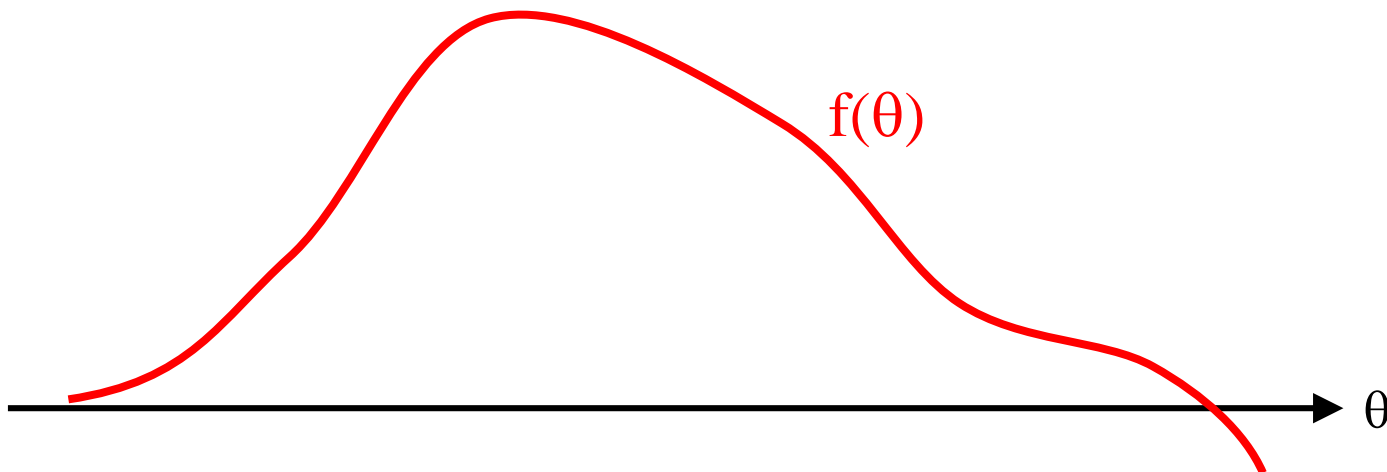
Hard to  
maximize



# Minorization-Maximization Algorithm

# Purpose of Minorization- Maximization (MM) Algorithm

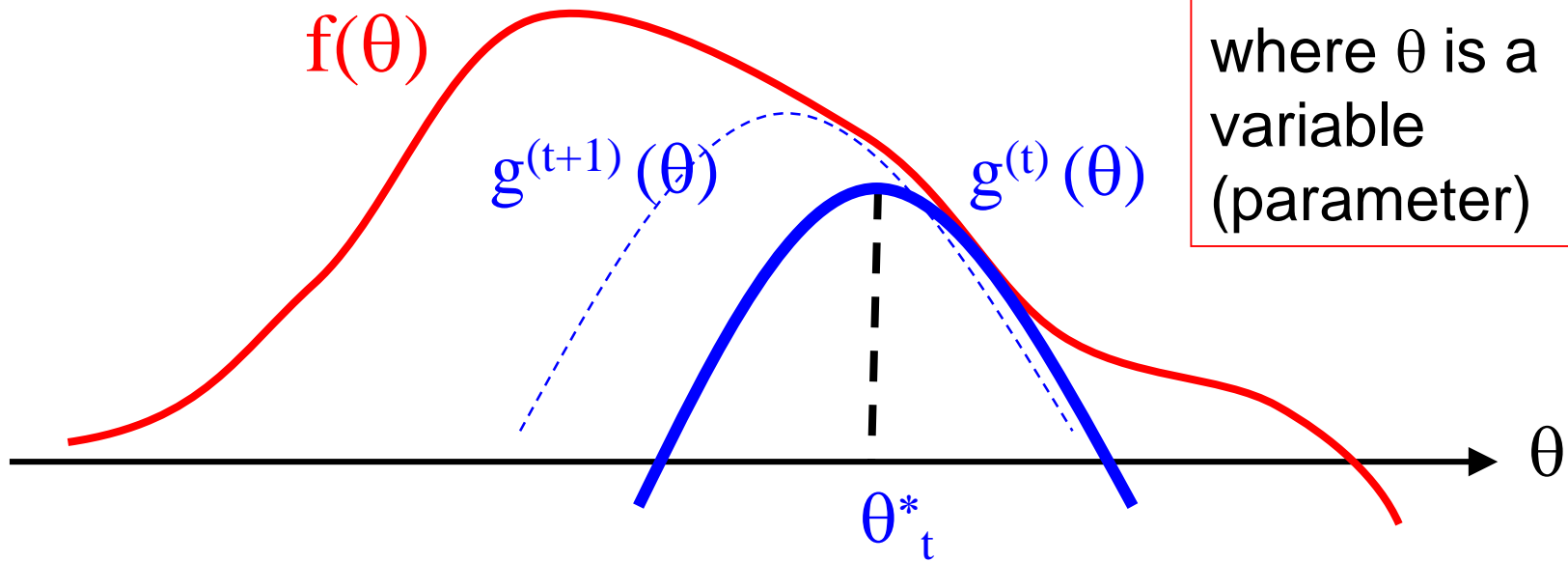
- Assume you want to maximize a function  $f(\theta)$  with respect to  $\theta$
- But it is hard to maximize  $f(\theta)$  directly..
  - because it is difficult to compute the 1<sup>st</sup> and 2<sup>nd</sup> derivatives
  - because  $\theta$  is very high-dimensional
  - Newton's method cannot be applied





# Basic Idea of MM Algorithm

- Instead, consider maximizing **a sequence of easier (surrogate) functions**  $\{ g^{(t)}(\theta) \}$  ( $t=1,2,\dots$ )
  - For example, **quadratic** functions are optimized easily
- But, what conditions should  $\{ g^{(t)}(\theta) \}$  satisfy ?
- Obviously,  $g^{(t)}(\theta)$  depends on the current solution  $\Rightarrow g^{(t)}(\theta) = g(\theta | \theta^t)$



$\theta^{(t)}$  is a constant  
where  $\theta$  is a  
variable  
(parameter)

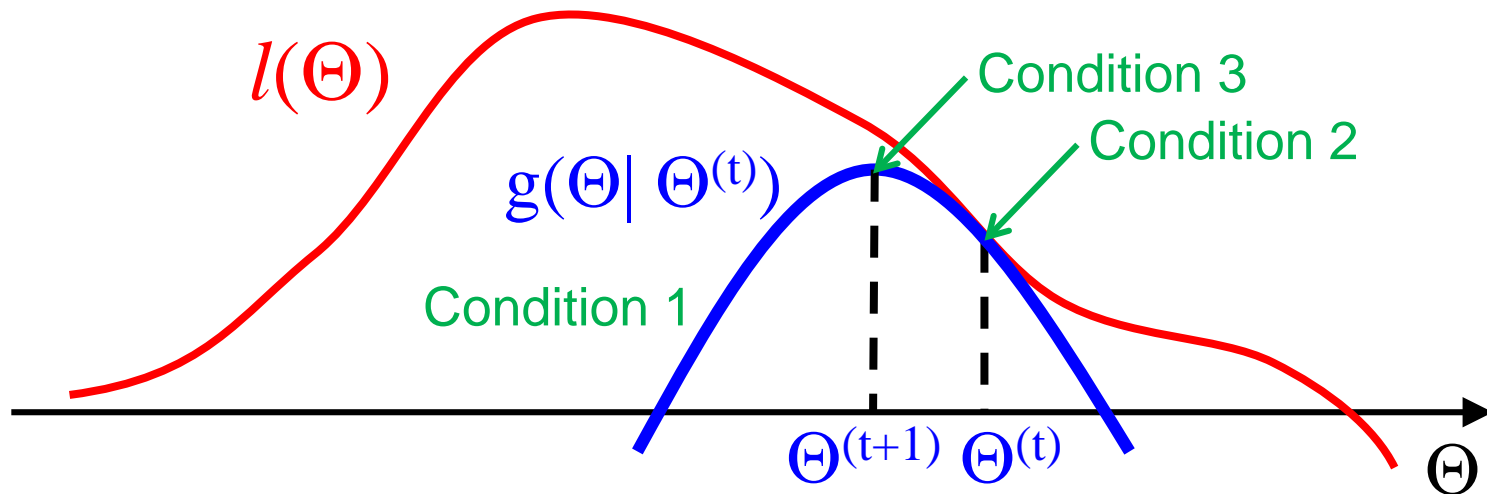
# Conditions for Surrogate Functions

[Condition 1]  $g(\theta | \theta^{(t)})$  must be a **lower bound** of  $f(\theta)$

i.e.,  $f(\theta) \geq g(\theta | \theta^{(t)})$  **for any  $\theta$**

[Condition 2]  $g(\theta^{(t)} | \theta^{(t)}) = f(\theta^{(t)})$

[Condition 3]  $g(\Theta | \Theta^{(t)})$  **can be easily maximized**



# Minorization-Maximization method

- Assume these 3 conditions are satisfied, then determine  $\Theta^{(t+1)}$  as,

$$\Theta^{(t+1)} = \arg \max_{\Theta} g(\Theta | \Theta^{(t)})$$

- Then the following inequality holds,

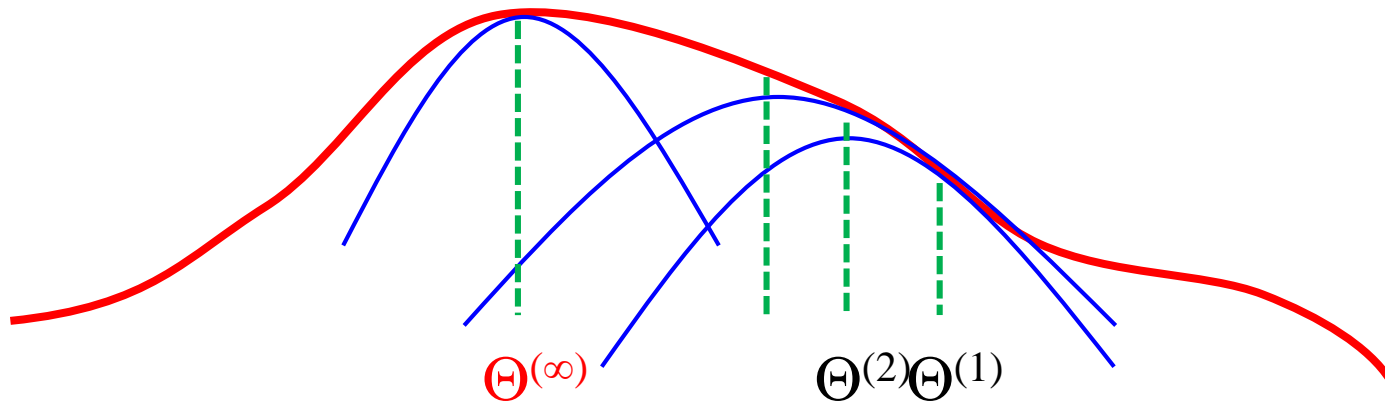
$$l(\Theta^{(t+1)}) \geq g(\Theta^{(t+1)} | \Theta^{(t)}) \geq g(\Theta^{(t)} | \Theta^{(t)}) = l(\Theta^{(t)})$$

↑  
Condition 1

↑  
 $\Theta^{(t+1)} = \arg \max_{\Theta} g(\Theta | \Theta^{(t)})$

↑  
Condition 2

- This leads to  $l(\Theta^{(1)}) \leq l(\Theta^{(2)}) \leq \dots \leq l(\Theta^{(\infty)})$  (Local) Maximum



# Jensen's inequality

# AM-GM Inequality

## (相加相乗平均の不等式)

- $x_1, x_2, \dots, x_n$  are non-negative real numbers
- AM : Arithmetic Mean  $\frac{1}{n}(x_1 + x_2 + \dots + x_n)$
- GM : Geometric Mean  $\sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$
- AM is greater than or equals to GM

$$\frac{1}{n}(x_1 + x_2 + \dots + x_n) \geq \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$$

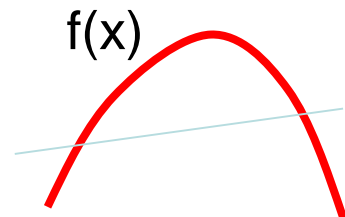
- Equality holds if and only if  $x_1 = x_2 = \dots = x_n$
- Taking the logarithm of both sides,

$$\log \frac{1}{n} \sum_{i=1}^n x_i \geq \frac{1}{n} \sum_{i=1}^n \log x_i \xrightarrow{\text{Assume } x \text{ is a r.v.}} \log E[x] \geq E[\log x]$$

log-sum
sum-log

# Jensen's Inequality

- Assume  $X$  is a random variable and  $f(X)$  is a **concave function**, then Jensen's inequality below holds



$$f(E[X]) \geq E[f(X)]$$

- If  $f(X)$  is **strictly concave**, equality holds if and only if  $X = \text{const.}$
- Now consider  $X$  is a **discrete** r.v. and  $p(X=x_i) = q_i$  and  $f(X) = \log X$  ( $\log$  is strictly concave)

$$f(E[x]) = \log\left(\sum_{i=1}^M q_i \cdot x_i\right) \geq E[f(x)] = \sum_{i=1}^M q_i \cdot \log x_i$$

Equality holds when  $x_1 = x_2 = \dots = x_M$

# EM algorithm

# Apply Jensen's Inequality to Log-likelihood Function

- Get back to the maximization of log-likelihood function of **unsupervised** learning problem
  - Consider  $X$  is discrete r.v. to be simple

$$l(\Theta) \equiv \ln p(Y | \Theta) = \ln \sum_X p(Y, X | \Theta)$$

- Consider  $q(X)$  is a probabilistic distribution of  $X$

$$l(\Theta) = \ln \sum_X p(Y, X | \Theta) = \ln \sum_X q(X) \frac{p(Y, X | \Theta)}{q(X)}$$

$$= \ln \mathbb{E}_{q(X)} \left[ \frac{p(Y, X | \Theta)}{q(X)} \right]$$


$$\geq \mathbb{E}_{q(X)} \left[ \ln \frac{p(Y, X | \Theta)}{q(X)} \right]$$

Jensen's Inequality



# Equality Condition and Lower Bound

Equality holds when  $\frac{p(Y, X | \Theta)}{q(X)} = \text{const.}$


$$q(X) \propto p(Y, X | \Theta) = p(X | Y, \Theta) \quad \left( \because \sum_X q(X) = 1 \right)$$

Now, define the lower bound function  $g(\Theta | \Theta^{(t)}) :$

$$g(\Theta | \Theta^{(t)}) \equiv \mathbb{E}_{p(X|Y, \Theta^{(t)})} \left[ \ln \frac{p(Y, X | \Theta)}{p(X | Y, \Theta^{(t)})} \right]$$

[Condition 1]  $l(\Theta) \geq g(\Theta | \Theta^{(t)})$  for all  $\Theta \Rightarrow$  Satisfied

[Condition 2]  $g(\Theta^{(t)} | \Theta^{(t)}) = l(\Theta^{(t)}) \Rightarrow$  Satisfied

[Condition 3]  $g(\Theta | \Theta^{(t)})$  can be easily maximized

# Maximizing Lower Bound

Log-likelihood      Lower bound

$$l(\Theta) \geq g(\Theta | \Theta^{(t)}) \Rightarrow \text{Maximize w.r.t. } \Theta$$

$$\begin{aligned}
 g(\Theta | \Theta^{(t)}) &= \mathbb{E}_{p(X|Y, \Theta^{(t)})} \left[ \ln \frac{p(Y, X | \Theta)}{p(X | Y, \Theta^{(t)})} \right] && \text{Called "complete" log likelihood} \\
 &= \underbrace{\mathbb{E}_{p(X|Y, \Theta^{(t)})} [\ln p(Y, X | \Theta)]}_{\text{Function of } \Theta} - \underbrace{\mathbb{E}_{p(X|Y, \Theta^{(t)})} [\ln p(X | Y, \Theta^{(t)})]}_{\text{Constant w.r.t. } \Theta}
 \end{aligned}$$

$$\begin{array}{c}
 \text{|||} \\
 Q(\Theta | \Theta^{(t)})
 \end{array}
 \text{ Expected complete log-likelihood}$$

In practice, we maximize  $Q(\Theta | \Theta^{(t)})$  at each iteration

$$\Theta^{(t+1)} \leftarrow \arg \max_{\Theta} Q(\Theta | \Theta^{(t)})$$

# Summary: EM Algorithm in General

Given:

- Data :  $Y$
- Initial parameter values:  $\Theta^{(0)}$

Repeat until convergence

1. [E-step] Compute posterior dist.  $q^*(X) = p(X | Y, \Theta^{(t)})$

$$\text{and } Q(\Theta | \Theta^{(t)}) = E_{q^*(X)} [\ln p(Y, X | \Theta)]$$

2. [M-step] Maximize  $Q(\Theta | \Theta^{(t)})$  w.r.t.  $\Theta$  Expected complete log-likelihood

$$\Theta^{(t+1)} \leftarrow \arg \max_{\Theta} Q(\Theta | \Theta^{(t)})$$

3.  $t \leftarrow t + 1$

# EM Algorithm for LDS

# Problem of Learning LDS

- Given :
  - Observation sequence :  $y_{1:T}$
- Find :
  - System matrices :  $A, B$
  - Noise covariance matrices:  $Q, R$
  - Initial state distribution:  $p(x_1) = N(x_1 | m_0, V_0)$
  - State sequence:  $x_{1:T}$ 
    - Mean vectors and covariance matrices :  $m_{1:T}$  and  $V_{1:T}$

Linear (Gaussian)  
Dynamical System

$$x_t = Ax_{t-1} + w_t$$

$$y_t = Cx_t + v_t$$

$$w_t \sim N(0, Q)$$

$$v_t \sim N(0, R)$$

# Likelihood Function of LDS

- Model parameters :  $\Theta = \{A, B, Q, R, m_0, V_0\}$
- Data :  $D = \{y_{1:T}\}$
- Log-likelihood function:

$$l(\Theta | D) \equiv \ln p(y_{1:T} | \Theta) = \ln \int p(y_{1:T}, x_{1:T} | \Theta) dx_{1:T}$$

$$= \ln \int \underbrace{p(y_{1:T} | x_{1:T}, \Theta)} p(x_{1:T} | \Theta) dx_{1:T}$$

Integral in log !



Hard to optimize (maximize)

# Much Easier Version

## If we could observe state sequence

- Given :
  - Observation sequence :  $\mathbf{y}_{1:T}$
  - State sequence:  $\mathbf{x}_{1:T}$ 
    - Mean vectors and covariance matrices :  $\mathbf{m}_{1:T}$  and  $\mathbf{V}_{1:T}$
- Find :
  - System matrices :  $\mathbf{A}, \mathbf{B}$
  - Noise covariance matrices:  $\mathbf{Q}, \mathbf{R}$
  - Initial state distribution:  $p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 \mid \mathbf{m}_0, \mathbf{V}_0)$

# Likelihood Function of Easier Version

- Model parameters :  $\Theta = \{A, B, Q, R, m_0, V_0\}$
- Data :  $D = \{y_{1:T}, x_{1:T}\}$
- Log-likelihood function:

$$\begin{aligned} l(\Theta | D) &\overset{\text{complete data log-likelihood}}{=} \ln p(y_{1:T}, x_{1:T} | \Theta) = \ln p(y_{1:T} | x_{1:T}, \Theta) p(x_{1:T} | \Theta) \\ &= \ln \left( \prod_{t=1}^T p(y_t | x_t, \Theta) \cdot p(x_1 | \Theta) \prod_{t=1}^{T-1} p(x_{t+1} | x_t, \Theta) \right) \\ &= \ln p(x_1 | \Theta) + \sum_{t=1}^{T-1} \ln p(x_{t+1} | x_t, \Theta) + \sum_{t=1}^T \ln p(y_t | x_t, \Theta) \end{aligned}$$



# Likelihood Function of Easier Version (cont.)

$$\begin{aligned}l(\Theta \mid D) &= \ln p(\mathbf{x}_1 \mid \Theta) + \sum_{t=1}^{T-1} \ln p(\mathbf{x}_{t+1} \mid \mathbf{x}_t, \Theta) + \sum_{t=1}^T \ln p(\mathbf{y}_t \mid \mathbf{x}_t, \Theta) \\&= \ln N(\mathbf{x}_1 \mid \mathbf{m}_0, \mathbf{V}_0) + \sum_{t=1}^{T-1} \ln N(\mathbf{x}_{t+1} \mid \mathbf{A}\mathbf{x}_t, \mathbf{Q}) + \sum_{t=1}^T \ln N(\mathbf{y}_t \mid \mathbf{C}\mathbf{x}_t, \mathbf{R}) \\&= -\frac{1}{2} \left\{ (\mathbf{x}_1 - \mathbf{m}_0)^T \mathbf{V}_0^{-1} (\mathbf{x}_1 - \mathbf{m}_0) + \sum_{t=1}^{T-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^T \mathbf{Q}^{-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) \right. \\&\quad \left. + \sum_{t=1}^T (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t)^T \mathbf{R}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t) + \ln |\mathbf{V}_0| + (T-1) \ln |\mathbf{Q}| + T \ln |\mathbf{R}| \right\} \\&\quad + \text{const.}\end{aligned}$$

# MLE Solution for Easier Version

Setting derivatives of  $l(\Theta|D)$  w.r.t. parameters to 0

$$\frac{\partial l(\Theta)}{\partial \mathbf{m}_0} = \mathbf{V}_0^{-1}(\mathbf{x}_1 - \mathbf{m}_0) = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{m}}_0 = \mathbf{x}_1$$

$$\frac{\partial l(\Theta)}{\partial \mathbf{V}_0} = -\frac{1}{2} \mathbf{V}_0^{-1} \left\{ (\mathbf{x}_1 - \mathbf{m}_0)(\mathbf{x}_1 - \mathbf{m}_0)^T \mathbf{V}_0^{-1} - \mathbf{I} \right\} = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{V}}_0 = \mathbf{0}$$

$$\frac{\partial l(\Theta)}{\partial \mathbf{A}} = \mathbf{Q}^{-1} \sum_{t=1}^{T-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) \mathbf{x}_t^T = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{A}} = \left( \sum_{t=1}^{T-1} \mathbf{x}_{t+1} \mathbf{x}_t^T \right) \left( \sum_{t=1}^{T-1} \mathbf{x}_t \mathbf{x}_t^T \right)^{-1}$$

$$\frac{\partial l(\Theta)}{\partial \mathbf{Q}} = -\frac{1}{2} \mathbf{Q}^{-1} \left\{ \sum_{t=1}^{T-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^T \mathbf{Q}^{-1} - (T-1)\mathbf{I} \right\} = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{Q}} = \frac{1}{T-1} \sum_{t=1}^{T-1} (\mathbf{x}_{t+1} - \hat{\mathbf{A}}\mathbf{x}_t)(\mathbf{x}_{t+1} - \hat{\mathbf{A}}\mathbf{x}_t)^T$$

$$\frac{\partial l(\Theta)}{\partial \mathbf{C}} = \mathbf{R}^{-1} \sum_{t=1}^T (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t) \mathbf{x}_t^T = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{C}} = \left( \sum_{t=1}^T \mathbf{y}_t \mathbf{x}_t^T \right) \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^T \right)^{-1}$$

$$\frac{\partial l(\Theta)}{\partial \mathbf{R}} = -\frac{1}{2} \mathbf{R}^{-1} \left\{ \sum_{t=1}^T (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t)(\mathbf{y}_t - \mathbf{C}\mathbf{x}_t)^T \mathbf{R}^{-1} - T \cdot \mathbf{I} \right\} = \mathbf{0} \quad \Rightarrow \quad \hat{\mathbf{R}} = \frac{1}{T} \sum_{t=1}^T (\mathbf{y}_t - \hat{\mathbf{C}}\mathbf{x}_t)(\mathbf{y}_t - \hat{\mathbf{C}}\mathbf{x}_t)^T$$

Can be solved analytically

# (Review) EM Algorithm in General

Given:

Observation sequence

- Data :  $Y = \mathbf{y}_{1:T}$
- Initial parameter values:  $\Theta^{(0)} = \{\mathbf{A}^{(0)}, \mathbf{B}^{(0)}, \mathbf{Q}^{(0)}, \mathbf{R}^{(0)}, \mathbf{m}_0^{(0)}, \mathbf{V}_0^{(0)}\}$

Repeat until convergence

Smoothed dist.

1. [E-step] Compute posterior dist.  $q^*(X) = p(X | Y, \Theta^{(t)})$

$$\text{and } Q(\Theta | \Theta^{(t)}) = E_{q^*(X)} [\ln p(Y, X | \Theta)]$$

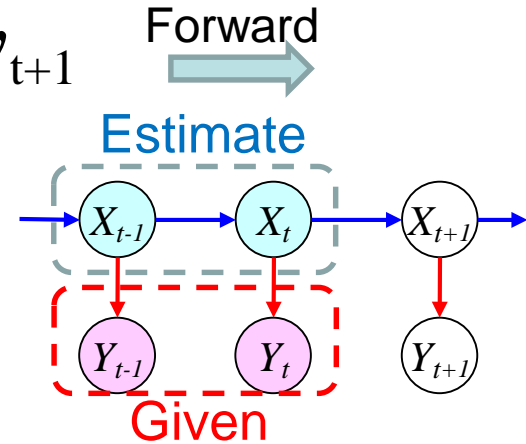
2. [M-step] Maximize  $Q(\Theta | \Theta^{(t)})$  w.r.t.  $\Theta$

$$\Theta^{(t+1)} \leftarrow \arg \max_{\Theta} Q(\Theta | \Theta^{(t)})$$

3.  $t \leftarrow t + 1$

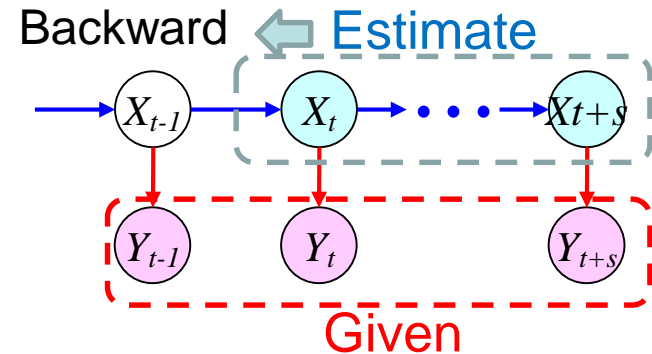
# E-step(1): Linear Kalman Filter

- Input :  $p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t)$  and  $\mathbf{y}_{t+1}$
- Prediction:  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{m}_t, \mathbf{P}_t)$ 
  - where  $\mathbf{P}_t = \mathbf{A}\mathbf{V}_t\mathbf{A}^T + \mathbf{Q}$
- Update :
 
$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{m}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{C}\mathbf{A}\mathbf{m}_t), (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{C})\mathbf{P}_t)$$
  - where  $\mathbf{K}_{t+1} = \mathbf{P}_t\mathbf{C}^T(\mathbf{C}\mathbf{P}_t\mathbf{C}^T + \mathbf{R})^{-1}$
- Output :  $p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{m}_{t+1}, \mathbf{V}_{t+1})$ 
  - where
 
$$\begin{cases} \mathbf{m}_{t+1} = \mathbf{A}\mathbf{m}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{C}\mathbf{A}\mathbf{m}_t) \\ \mathbf{V}_{t+1} = (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{C})\mathbf{P}_t \end{cases}$$



# E-step(2) : RTS Smoothing

- Assume filtered dist.  
 $p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t)$   
 have been computed  
 by Kalman filtering



- At terminal time,  $p(\mathbf{x}_T | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_T | \mathbf{m}_T, \mathbf{V}_T) = \mathcal{N}(\mathbf{x}_T | \hat{\mathbf{m}}_T, \hat{\mathbf{V}}_T)$
- For  $t=T-1$  to 1, repeat the following computation

$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t | \hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t)$$

where

$$\begin{cases} \hat{\mathbf{m}}_t = \mathbf{m}_t + \mathbf{J}_t (\hat{\mathbf{m}}_{t+1} - \mathbf{A} \mathbf{m}_t) \\ \hat{\mathbf{V}}_t = \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t \mathbf{A} \mathbf{V} = \mathbf{V}_t + \mathbf{J}_t (\hat{\mathbf{V}}_{t+1} - \mathbf{P}_t) \mathbf{J}_t^T \end{cases}$$

and covariance between  $\mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  given  $\mathbf{y}_{1:T}$ :

$$\text{cov}[\mathbf{x}_t, \mathbf{x}_{t+1}] = \mathbf{J}_t \hat{\mathbf{V}}_{t+1} \Rightarrow \text{We require this later}$$

# E-step (3) : Miscellaneous Expectations

In M-step, we will require the following expectations:


$$\left\{ \begin{array}{l} \mathbb{E}_{q^{(t)}(\mathbf{x}_{1:T})}[\mathbf{x}_t] = \hat{\mathbf{m}}_t \\ \mathbb{E}_{q^{(t)}(\mathbf{x}_{1:T})}[\mathbf{x}_t \mathbf{x}_t^T] = \text{var}(\mathbf{x}_t) + \hat{\mathbf{m}}_t \hat{\mathbf{m}}_t^T = \hat{\mathbf{V}}_t + \hat{\mathbf{m}}_t \hat{\mathbf{m}}_t^T \\ \mathbb{E}_{q^{(t)}(\mathbf{x}_{1:T})}[\mathbf{x}_t \mathbf{x}_{t+1}^T] = \text{cov}[\mathbf{x}_t, \mathbf{x}_{t+1}] + \hat{\mathbf{m}}_t \hat{\mathbf{m}}_{t+1}^T = \mathbf{J}_t \hat{\mathbf{V}}_{t+1} + \hat{\mathbf{m}}_t \hat{\mathbf{m}}_{t+1}^T \end{array} \right.$$

where,  $q^{(t)}(\mathbf{x}_{1:T}) = p(\mathbf{x}_{1:T} \mid \mathbf{y}_{1:T}, \Theta^{(t)})$

is the joint smoothed distribution of state sequence,  
given observation sequence and estimated parameter  
at t-th iteration

# M-step(1) : Complete Log-Likelihood

We've got already the complete log-likelihood

$$\begin{aligned} \ln p(\mathbf{y}_{1:T}, \mathbf{x}_{1:T} \mid \Theta) \\ = -\frac{1}{2} \left\{ (\mathbf{x}_1 - \mathbf{m}_0)^T \mathbf{V}_0^{-1} (\mathbf{x}_1 - \mathbf{m}_0) + \sum_{t=1}^{T-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^T \mathbf{Q}^{-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) \right. \\ \left. + \sum_{t=1}^T (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t)^T \mathbf{R}^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t) + \ln |\mathbf{V}_0| + (T-1) \ln |\mathbf{Q}| + T \ln |\mathbf{R}| \right\} \\ + \text{const.} \end{aligned}$$


$$Q(\Theta \mid \Theta^{(t)}) = \mathbb{E}_{q^*(X)} [\ln p(\mathbf{y}_{1:T}, \mathbf{x}_{1:T} \mid \Theta)]$$

Compute the expectation of this function w.r.t. state sequence  $\{\mathbf{x}_{1:T}\}$  and maximize w.r.t. parameters  $\Theta = \{\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R}, \mathbf{m}_0, \mathbf{V}_0\}$  !

# M-step (2) : Update of Parameters

$$\frac{\partial Q}{\partial \mathbf{m}_0} = \mathbf{V}_0^{-1} (\mathbb{E}[\mathbf{x}_1] - \mathbf{m}_0) = \mathbf{0} \quad \Rightarrow \quad \mathbf{m}_0^{(t+1)} = \mathbb{E}[\mathbf{x}_1] = \hat{\mathbf{m}}_1$$

$$\frac{\partial Q}{\partial \mathbf{V}_0} = -\frac{1}{2} \mathbf{V}_0^{-1} \left\{ \mathbb{E}[(\mathbf{x}_1 - \mathbf{m}_0)(\mathbf{x}_1 - \mathbf{m}_0)^T] \mathbf{V}_0^{-1} - \mathbf{I} \right\} = \mathbf{0} \quad \Rightarrow \quad \mathbf{V}_0^{(t+1)} = \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^T] - \mathbf{m}_0 \mathbf{m}_0^T = \hat{\mathbf{V}}_1$$

$$\frac{\partial Q}{\partial \mathbf{A}} = -\frac{1}{2} \mathbf{Q}^{-1} \sum_{t=1}^{T-1} (\mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^T] - \mathbf{A} \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T]) = \mathbf{0} \quad \Rightarrow \quad \mathbf{A}^{(t+1)} = \left( \sum_{t=1}^{T-1} \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^T] \right) \left( \sum_{t=1}^{T-1} \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T] \right)^{-1}$$

$$\frac{\partial Q}{\partial \mathbf{Q}} = -\frac{1}{2} \mathbf{Q}^{-1} \left\{ \sum_{t=1}^{T-1} \mathbb{E}[(\mathbf{x}_{t+1} - \mathbf{A} \mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{A} \mathbf{x}_t)^T] \mathbf{Q}^{-1} - (T-1) \mathbf{I} \right\} = \mathbf{0}$$

$$\Rightarrow \quad \mathbf{Q}^{(t+1)} = \frac{1}{T-1} \sum_{t=1}^{T-1} \left\{ \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_{t+1}^T] - \mathbf{A}^{(t+1)} \mathbb{E}[\mathbf{x}_t \mathbf{x}_{t+1}^T] - \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^T] \mathbf{A}^{(t+1)T} + \mathbf{A}^{(t+1)} \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T] \mathbf{A}^{(t+1)T} \right\}$$

$$\frac{\partial Q}{\partial \mathbf{C}} = \mathbf{R}^{-1} \sum_{t=1}^T (\mathbf{y}_t \mathbb{E}[\mathbf{x}_t^T] - \mathbf{C} \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T]) = \mathbf{0} \quad \Rightarrow \quad \mathbf{C}^{(t+1)} = \left( \sum_{t=1}^T \mathbf{y}_t \mathbb{E}[\mathbf{x}_t^T] \right) \left( \sum_{t=1}^T \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T] \right)^{-1}$$

$$\frac{\partial Q}{\partial \mathbf{R}} = -\frac{1}{2} \mathbf{R}^{-1} \left\{ \sum_{t=1}^T \mathbb{E}[(\mathbf{y}_t - \mathbf{C} \mathbf{x}_t)(\mathbf{y}_t - \mathbf{C} \mathbf{x}_t)^T] \mathbf{R}^{-1} - T \cdot \mathbf{I} \right\} = \mathbf{0}$$

$$\Rightarrow \quad \mathbf{R}^{(t+1)} = \frac{1}{T} \sum_{t=1}^T \left\{ \mathbf{y}_t \mathbf{y}_t^T - \mathbf{C}^{(t+1)} \mathbb{E}[\mathbf{x}_t] \mathbf{y}_t^T - \mathbf{y}_t \mathbb{E}[\mathbf{x}_t^T] \mathbf{C}^{(t+1)T} + \mathbf{C}^{(t+1)} \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^T] \mathbf{C}^{(t+1)T} \right\} \quad 64$$



# Summary

- EM algorithm for learning linear dynamical systems was derived
- E-step (inference) is actually Kalman (RTS) smoother
- M-step is similar to ordinary maximum likelihood estimation for supervised learning
- This framework can be extended to many complicated (i.e., non-linear, switching, etc.) models
  - Often some approximation is necessary, though.
  - We will see such extensions next time