# Sequential Data Modeling

## " Linear Dynamical Systems"
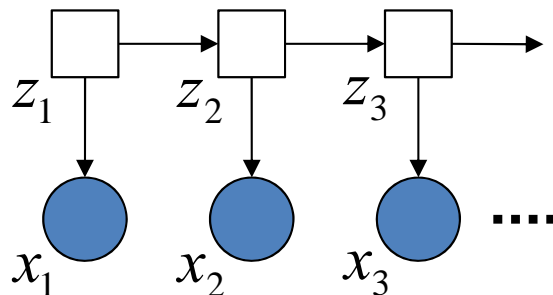
# Basic Techniques



Discrete latent variables

Continuous latent variables

Markov model

Mixture model (*e.g.*, GMM)

$z_1$ $z_2$ $z_3$

$x_1$ $x_2$ $x_3$ ....

Factor analysis (FA)

$z_1$ $z_2$ $z_3$

$x_1$ $x_2$ $x_3$ ....

hidden Markov model (HMM)

$z_1$ $z_2$ $z_3$

$x_1$ $x_2$ $x_3$ ....

Linear dynamical systems (LDS)

$z_1$ $z_2$ $z_3$

$x_1$ $x_2$ $x_3$ ....

# Linear Dynamical Systems (LDS)

# Basic Techniques

# Assume Unobservable Data Sequence

- Extraction of an underlying data sequence from observable data sequence suffering from noise

How to model this?

$x_t$

Observation data sequence
**w/ noise**

$t$

How to extract this?

$z_t$

Underlying data sequence
(**unobservable**)

$t$

# How to Model Sequential Data?

- To model sequential data...

  - Need to consider sample order

  - Need to model a very high-dimensional space of joint data over a sequence

  - Need to deal with various lengths of sequential data



Factor analysis (FA)

$z_1$ $z_2$ $z_3$

$x_1$ $x_2$ $x_3$ ....

➢Lots of dependencies...
➢Length varies...

$x_1$ $x_2$ $x_3$ ....

$$p(\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_T \mid \boldsymbol{\lambda}) = \prod_{t=1}^{T} p(\boldsymbol{x}_t \mid \boldsymbol{\lambda})$$

$$p(\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_T \mid \boldsymbol{\lambda}) = \ ?$$

How to model this $p.d.f.$?

# Linear Dynamical Systems

- **Markov process** to model a sequence of continuous latent variables
- **Linear equation** to model state transition and mapping from a state space into an observation space



**State space**

$$z_t = Az_{t-1} + n_t$$

Transition noise:
$$n_t \sim \mathcal{N}(n_t; 0, \Gamma)$$

State transition *p.d.f.*
$$p(z_t \mid z_{t-1}) = \mathcal{N}(z_t; Az_{t-1}, \Gamma)$$

**Observation space**

$$x_t = Wz_t + e_t$$

Observation noise:
$$e_t \sim \mathcal{N}(e_t; 0, \Sigma)$$

Emission *p.d.f.*
$$p(x_t \mid z_t) = \mathcal{N}(x_t; Wz_t, \Sigma)$$

# *p.d.f.*s in Linear Dynamical Systems

- Sequence of observation data : $x_{1:T} = \{x_1, x_2, \cdots, x_T\}$
- Sequence of latent variables : $z_{1:T} = \{z_1, z_2, \cdots, z_T\}$
- *p.d.f.* of observation data :

$$p(x_{1:T}) = \int p(x_{1:T} \mid z_{1:T}) p(z_{1:T}) \, dz_{1:T}$$

Marginalization over a sequence of latent variables

$$= \int \left[ \prod_{t=1}^{T} p(x_t \mid z_t) \right] \left[ p(z_1) \prod_{t=2}^{T} p(z_t \mid z_{t-1}) \right] dz_{1:T}$$

Emission *p.d.f.*s:
$$p(x_t \mid z_t) = \mathcal{N}(x_t; Wz_t, \Sigma)$$

Transition *p.d.f.*s:
$$\begin{cases} p(z_t \mid z_{t-1}) = \mathcal{N}(z_t; Az_{t-1}, \Gamma) & t \geq 2 \\ p(z_1) = \mathcal{N}(z_1; \mu_0, P_0) & t = 1 \end{cases}$$

# Kalman Filtering

Propagate uncertainty
from past to future

# How to Recursively Calculate Likelihood?

- Likelihood function for the observation data sequence:

$$p(\boldsymbol{x}) = \int p(\boldsymbol{x} \mid \boldsymbol{z})p(\boldsymbol{z})\,\mathrm{d}\boldsymbol{z} = \mathcal{N}\left(\boldsymbol{x}; \tilde{\boldsymbol{W}}\tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{z}}_0, \tilde{\boldsymbol{W}}\tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{\Gamma}}\tilde{\boldsymbol{A}}^{-\top}\tilde{\boldsymbol{W}}^{\top} + \tilde{\boldsymbol{\Sigma}}\right)$$

Batch-type calculation (w/ all data over an sequence) is assumed but **frame-by-frame calculation** will be required in some applications, such as real-time signal processing…

- Original form of the likelihood function:

$$p(\boldsymbol{x}_{1:T}) = \int p(\boldsymbol{x}_{1:T} \mid \boldsymbol{z}_{1:T})p(\boldsymbol{z}_{1:T})\,\mathrm{d}\boldsymbol{z}_{1:T}$$

$$= \int \left[\prod_{t=1}^{T} p(\boldsymbol{x}_t \mid \boldsymbol{z}_t)\right]\left[p(\boldsymbol{z}_1)\prod_{t=2}^{T} p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1})\right]\mathrm{d}\boldsymbol{z}_{1:T}$$



11

# Kalman Filtering (Forward Algorithm)

- Posterior $p.d.f.$s given all past observation data ($i.e.$, $p(z_t \mid x_{1:t}) = \alpha(z_t)$) determined in Kalman filtering



Kalman filtering

# Prediction and Update

- **Prediction step**
  - Predict distribution of latent variables at frame $t$ from all past observation data



- **Update step**
  - Update distribution of latent variables at frame $t$ using current observation data as well as all past observation data

# Predicted and Updated $p.d.f.$s

- **Predicted** $p.d.f.$

$$p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1}) = \int p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}) p(\boldsymbol{z}_{t-1} \mid \boldsymbol{x}_{1:t-1}) \mathrm{d}\boldsymbol{z}_t = \mathcal{N}\left(\boldsymbol{z}_t; \boldsymbol{\mu}_{t|t-1}, \boldsymbol{P}_{t|t-1}\right)$$

Predicted mean : $\boldsymbol{\mu}_{t|t-1} = \boldsymbol{A}\boldsymbol{\mu}_{t-1}$

Predicted covariance : $\boldsymbol{P}_{t|t-1} = \boldsymbol{A}\boldsymbol{P}_{t-1}\boldsymbol{A}^\top + \boldsymbol{\Gamma}$

- **Updated** $p.d.f.$

**Posterior $\propto$ Likelihood x Prior**

$$\alpha(\boldsymbol{z}_t) = p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t}) = \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{\mu}_t, \boldsymbol{P}_t) \propto p(\boldsymbol{x}_t \mid \boldsymbol{z}_t) p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1})$$

Kalman gain matrix : $\boldsymbol{K}_t = \boldsymbol{P}_{t|t-1}\boldsymbol{W}^\top\left(\boldsymbol{W}\boldsymbol{P}_{t|t-1}\boldsymbol{W}^\top + \boldsymbol{\Sigma}\right)^{-1}$

Updated mean : $\boldsymbol{\mu}_t = \boldsymbol{\mu}_{t|t-1} + \boldsymbol{K}_t\left(\boldsymbol{x}_t - \boldsymbol{W}\boldsymbol{\mu}_{t|t-1}\right)$

Updated covariance : $\boldsymbol{P}_t = \left(\boldsymbol{I} - \boldsymbol{K}_t\boldsymbol{W}\right)\boldsymbol{P}_{t|t-1}$

**Error between predicted and observed data**

# Likelihood Calculation

- **Conditional** *p.d.f.* **of observation data**

$$p(\boldsymbol{x}_t \mid \boldsymbol{x}_{1:t-1}) = \int p(\boldsymbol{x}_t \mid \boldsymbol{z}_t) p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1}) \, \mathrm{d}\boldsymbol{z}_t = \mathcal{N}\left(\boldsymbol{x}_t ; \boldsymbol{W}\boldsymbol{\mu}_{t|t-1}, \boldsymbol{W}\boldsymbol{P}_{t|t-1}\boldsymbol{W}^{\top} + \boldsymbol{\Sigma}\right)$$

Emission *p.d.f.* :  $p(\boldsymbol{x}_t \mid \boldsymbol{z}_t) = \mathcal{N}\left(\boldsymbol{x}_t ; \boldsymbol{W}\boldsymbol{z}_t, \boldsymbol{\Sigma}\right)$

Predicted *p.d.f.* :  $p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1}) = \mathcal{N}\left(\boldsymbol{z}_t ; \boldsymbol{\mu}_{t|t-1}, \boldsymbol{P}_{t|t-1}\right)$

- **Recursive likelihood calculation**

$$p(\boldsymbol{x}_{1:T}) = \underbrace{\underbrace{p(\boldsymbol{x}_1) p(\boldsymbol{x}_2 \mid \boldsymbol{x}_1)}_{p(\boldsymbol{x}_1, \boldsymbol{x}_2)} p(\boldsymbol{x}_3 \mid \boldsymbol{x}_1, \boldsymbol{x}_2) \cdots p(\boldsymbol{x}_{T-1} \mid \boldsymbol{x}_{1:T-2})}_{p(\boldsymbol{x}_{1:T-1})} p(\boldsymbol{x}_T \mid \boldsymbol{x}_{1:T-1})$$

$$\underbrace{\qquad}_{p(\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{x}_3)} \cdots$$

# Kalman Smoothing

# How to Calculate Posterior w/ All Data?

- Posterior $p.d.f.$s given all past observation data determined in Kalman filtering



- How to calculate posterior $p.d.f.$s calculated w/ all observation data?

# Kalman Smoothing (Backward Algorithm)

- Calculate smoothed *p.d.f.* at frame $t$ using both smoothed *p.d.f.* at frame $t+1$ and updated *p.d.f.* at frame $t$ determined in Kalman filtering



See appendix

Kalman smoothing

$$\gamma(z_t) = \mathcal{N}\left(z_t; \hat{\mu}_t, \hat{P}_t\right)$$

$$J_t = P_t A^\mathsf{T} P_{t+1|t}^{-1}$$

Smoothed mean:
$$\hat{\mu}_t = \mu_t + J_t\left(\hat{\mu}_{t+1} - \mu_{t+1|t}\right)$$

Smoothed covariance :
$$\hat{P}_t = P_t + J_t\left(\hat{P}_{t+1} - P_{t+1|t}\right)J_t^\mathsf{T}$$

# Model Training

# EM Algorithm

- Likelihood function

$$p(\boldsymbol{x}_{1:T} \mid \lambda) = \int p(\boldsymbol{x}_{1:T}, \boldsymbol{z}_{1:T} \mid \lambda)\, \mathrm{d}\boldsymbol{z}_{1:T}$$

$$= \int \left[ \prod_{t=1}^{T} p(\boldsymbol{x}_t \mid \boldsymbol{z}_t, \{\boldsymbol{W}, \boldsymbol{\Sigma}\}) \right] \left[ p(\boldsymbol{z}_1 \mid \{\boldsymbol{\mu}_0, \boldsymbol{P}_0\}) \prod_{t=2}^{T} p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}, \{\boldsymbol{A}, \boldsymbol{\Gamma}\}) \right] \mathrm{d}\boldsymbol{z}_{1:T}$$

- Iterative maximization of **lower bound**

$$\ln p(\boldsymbol{x}_{1:T} \mid \lambda) = \ln \int p(\boldsymbol{x}_{1:T}, \boldsymbol{z}_{1:T} \mid \lambda)\, \mathrm{d}\boldsymbol{z}_{1:T}$$

$$\geq \int q(\boldsymbol{z}_{1:T}) \ln \frac{p(\boldsymbol{x}_{1:T}, \boldsymbol{z}_{1:T} \mid \lambda)}{q(\boldsymbol{z}_{1:T})}\, \mathrm{d}\boldsymbol{z}_{1:T} = \mathcal{L}(q, \lambda)$$

**E-step**: Set $q$ to the posterior *p.d.f.* calculated w/ current model parameters

$$\hat{q}(\boldsymbol{z}_{1:T}) = p(\boldsymbol{z}_{1:T} \mid \boldsymbol{x}_{1:T}, \lambda_{\text{old}})$$

**M-step**: Maximize auxiliary function with respect to model parameters

$$\hat{\lambda} = \arg\max_{\lambda} \int \hat{q}(\boldsymbol{z}_{1:T}) \ln\{p(\boldsymbol{x}_{1:T}, \boldsymbol{z}_{1:T} \mid \lambda)\}\, \mathrm{d}\boldsymbol{z}_{1:T}$$

# E-Step: Update $q$

- Calculation of posterior $p.d.f.$s using a model parameter set $\boldsymbol{\lambda}_{\text{old}}$

  - Posterior $p.d.f.$ of $\boldsymbol{z}_t$ calculated in <span style="color:red">Kalman smoothing</span>

$$\hat{q}(\boldsymbol{z}_t) = p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:T}, \boldsymbol{\lambda}_{\text{old}}) = \gamma(\boldsymbol{z}_t) = \mathcal{N}(\boldsymbol{z}_t; \hat{\boldsymbol{\mu}}_t, \hat{\boldsymbol{P}}_t)$$

  - Joint posterior $p.d.f.$ of $\boldsymbol{z}_{t-1}$ and $\boldsymbol{z}_t$ also calculated in <span style="color:red">Kalman smoothing</span>

$$\hat{q}(\boldsymbol{z}_{t-1}, \boldsymbol{z}_t) = p(\boldsymbol{z}_{t-1}, \boldsymbol{z}_t \mid \boldsymbol{x}_{1:T}, \boldsymbol{\lambda}_{\text{old}})$$

$$= \mathcal{N}\left( \begin{bmatrix} \boldsymbol{z}_{t-1} \\ \boldsymbol{z}_t \end{bmatrix}; \begin{bmatrix} \hat{\boldsymbol{\mu}}_{t-1} \\ \hat{\boldsymbol{\mu}}_t \end{bmatrix}, \begin{bmatrix} \hat{\boldsymbol{P}}_{t-1} & \boldsymbol{J}_{t-1}\hat{\boldsymbol{P}}_t \\ \hat{\boldsymbol{P}}_t \boldsymbol{J}_{t-1}^{\mathsf{T}} & \hat{\boldsymbol{P}}_t \end{bmatrix} \right)$$

Note that $\hat{q}(\boldsymbol{z}_t) = \int \hat{q}(\boldsymbol{z}_{t-1}, \boldsymbol{z}_t)\, \mathrm{d}\boldsymbol{z}_{t-1}$

# M-Step: Update $\lambda$

- Maximization of the following auxiliary function

$$Q(\lambda_{\text{old}}, \lambda) = \int \hat{q}(z_{1:T}) \left\{ \sum_{t=1}^{T} \underbrace{\ln p(x_t \mid z_t, \lambda)}_{\textcircled{1}} + \underbrace{\ln p(z_1 \mid \lambda)}_{\textcircled{2}} + \sum_{t=2}^{T} \underbrace{\ln p(z_t \mid z_{t-1}, \lambda)}_{\textcircled{3}} \right\} dz_{1:T}$$

$$= \int \hat{q}(z_{1:T}) \left\{ \sum_{t=1}^{T} \underbrace{\frac{1}{2}\ln|\Sigma^{-1}| - \frac{1}{2}(x_t - Wz_t)^{\top} \Sigma^{-1}(x_t - Wz_t)}_{\textcircled{1}} \right.$$

$$+ \underbrace{\frac{1}{2}\ln|P_0^{-1}| - \frac{1}{2}(z_1 - \mu_0)^{\top} P_0^{-1}(z_1 - \mu_0)}_{\textcircled{2}}$$

$$\left. + \sum_{t=2}^{T} \underbrace{\frac{1}{2}\ln|\Gamma^{-1}| - \frac{1}{2}(z_t - Az_{t-1})^{\top} \Gamma^{-1}(z_t - Az_{t-1})}_{\textcircled{3}} \right\} dz_{1:T}$$

# Expansion of Auxiliary Function

$$Q(\lambda_{\text{old}}, \lambda) = \frac{1}{2} \Big\{ T \ln|\boldsymbol{\Sigma}^{-1}| - \text{tr}\Big[ \boldsymbol{\Sigma}^{-1} \langle \boldsymbol{x}_t \boldsymbol{x}_t^\mathsf{T} \rangle_{1:T} + \boldsymbol{W}^\mathsf{T} \boldsymbol{\Sigma}^{-1} \boldsymbol{W} \langle \boldsymbol{z}_t \boldsymbol{z}_t^\mathsf{T} \rangle_{1:T}$$

$$- \boldsymbol{\Sigma}^{-1} \boldsymbol{W} \langle \langle \boldsymbol{z}_t \rangle \boldsymbol{x}_t^\mathsf{T} \rangle_{1:T} - \langle \langle \boldsymbol{z}_t \rangle \boldsymbol{x}_t^\mathsf{T} \rangle_{1:T}^\mathsf{T} \boldsymbol{W}^\mathsf{T} \boldsymbol{\Sigma}^{-1} \Big]$$

$$+ \ln|\boldsymbol{P}_0^{-1}| - \text{tr}\Big[ \boldsymbol{P}_0^{-1} \langle \boldsymbol{z}_1 \boldsymbol{z}_1^\mathsf{T} \rangle + \boldsymbol{P}_0^{-1} \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^\mathsf{T} \Big] + \langle \boldsymbol{z}_1 \rangle^\mathsf{T} \boldsymbol{P}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0^\mathsf{T} \boldsymbol{P}_0^{-1} \langle \boldsymbol{z}_1 \rangle$$

$$+ (T-1) \ln|\boldsymbol{\Gamma}^{-1}| - \text{tr}\Big[ \boldsymbol{\Gamma}^{-1} \langle \boldsymbol{z}_t \boldsymbol{z}_t^\mathsf{T} \rangle_{2:T} + \boldsymbol{A}^\mathsf{T} \boldsymbol{\Gamma}^{-1} \boldsymbol{A} \langle \boldsymbol{z}_{t-1} \boldsymbol{z}_{t-1}^\mathsf{T} \rangle_{2:T}$$

$$- \boldsymbol{\Gamma}^{-1} \boldsymbol{A} \langle \boldsymbol{z}_{t-1} \boldsymbol{z}_t^\mathsf{T} \rangle_{2:T} - \boldsymbol{A}^\mathsf{T} \boldsymbol{\Gamma}^{-1} \langle \boldsymbol{z}_{t-1} \boldsymbol{z}_t^\mathsf{T} \rangle_{2:T}^\mathsf{T} \Big] \Big\}$$

**Expectation:**

$$\langle \boldsymbol{z}_t \rangle = \hat{\boldsymbol{\mu}}_t$$

$$\langle \boldsymbol{z}_t \boldsymbol{z}_t^\mathsf{T} \rangle = \hat{\boldsymbol{P}}_t + \hat{\boldsymbol{\mu}}_t \hat{\boldsymbol{\mu}}_t^\mathsf{T}$$

$$\langle \boldsymbol{z}_{t-1} \boldsymbol{z}_t^\mathsf{T} \rangle = \boldsymbol{J}_{t-1} \hat{\boldsymbol{P}}_t + \hat{\boldsymbol{\mu}}_{t-1} \hat{\boldsymbol{\mu}}_t^\mathsf{T}$$

**Sufficient statistics:**

$$\langle \boldsymbol{z}_t \boldsymbol{z}_t^\mathsf{T} \rangle_{1:T} = \sum_{t=1}^{T} \langle \boldsymbol{z}_t \boldsymbol{z}_t^\mathsf{T} \rangle \qquad \langle \langle \boldsymbol{z}_t \rangle \boldsymbol{x}_t^\mathsf{T} \rangle_{1:T} = \sum_{t=1}^{T} \langle \boldsymbol{z}_t \rangle \boldsymbol{x}_t^\mathsf{T}$$

$$\langle \boldsymbol{z}_{t-1} \boldsymbol{z}_t^\mathsf{T} \rangle_{2:T} = \sum_{t=2}^{T} \langle \boldsymbol{z}_{t-1} \boldsymbol{z}_t^\mathsf{T} \rangle \qquad \langle \boldsymbol{x}_t \boldsymbol{x}_t^\mathsf{T} \rangle_{1:T} = \sum_{t=1}^{T} \boldsymbol{x}_t \boldsymbol{x}_t^\mathsf{T}$$

27

# ML Estimates of Model Parameters

**Initial parameters of transition *p.d.f.***

$$\hat{\boldsymbol{\mu}}_0 = \left\langle \boldsymbol{z}_1 \right\rangle$$

$$\hat{\boldsymbol{P}}_0 = \left\langle \boldsymbol{z}_1 \boldsymbol{z}_1^{\mathsf{T}} \right\rangle - \hat{\boldsymbol{\mu}}_0 \hat{\boldsymbol{\mu}}_0^{\mathsf{T}}$$

**Parameters of transition *p.d.f.***

$$\hat{\boldsymbol{A}} = \left\langle \boldsymbol{z}_t \boldsymbol{z}_{t-1}^{\mathsf{T}} \right\rangle_{2:T} \left\langle \boldsymbol{z}_{t-1} \boldsymbol{z}_{t-1}^{\mathsf{T}} \right\rangle_{2:T}^{-1}$$

$$\hat{\boldsymbol{\Gamma}} = \frac{1}{T-1} \left( \left\langle \boldsymbol{z}_t \boldsymbol{z}_t^{\mathsf{T}} \right\rangle_{2:T} + \hat{\boldsymbol{A}} \left\langle \boldsymbol{z}_{t-1} \boldsymbol{z}_{t-1}^{\mathsf{T}} \right\rangle_{2:T} \hat{\boldsymbol{A}}^{\mathsf{T}} - \hat{\boldsymbol{A}} \left\langle \boldsymbol{z}_{t-1} \boldsymbol{z}_t^{\mathsf{T}} \right\rangle_{2:T} - \left\langle \boldsymbol{z}_t \boldsymbol{z}_{t-1}^{\mathsf{T}} \right\rangle_{2:T} \hat{\boldsymbol{A}}^{\mathsf{T}} \right)$$

**Parameters of emission *p.d.f.***

$$\hat{\boldsymbol{W}} = \left\langle \boldsymbol{x}_t \left\langle \boldsymbol{z}_t \right\rangle^{\mathsf{T}} \right\rangle_{1:T} \left\langle \boldsymbol{z}_t \boldsymbol{z}_t^{\mathsf{T}} \right\rangle_{1:T}^{-1}$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{T} \left( \left\langle \boldsymbol{x}_t \boldsymbol{x}_t^{\mathsf{T}} \right\rangle_{1:T} + \hat{\boldsymbol{W}} \left\langle \boldsymbol{z}_t \boldsymbol{z}_t^{\mathsf{T}} \right\rangle_{1:T} \hat{\boldsymbol{W}}^{\mathsf{T}} - \hat{\boldsymbol{W}} \left\langle \left\langle \boldsymbol{z}_t \right\rangle \boldsymbol{x}_t^{\mathsf{T}} \right\rangle_{1:T} - \left\langle \boldsymbol{x}_t \left\langle \boldsymbol{z}_t \right\rangle^{\mathsf{T}} \right\rangle_{1:T} \hat{\boldsymbol{W}}^{\mathsf{T}} \right)$$

# Appendix

# Derivation of $p.d.f.$ of observation data

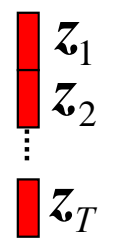# Emission $p.d.f.\ p(\pmb{x}_{1:\text{T}}|\pmb{z}_{1:\text{T}})$



- Vector form to represent a data sequence

Observation data:
$$\pmb{x} = \left[\pmb{x}_1^\top, \pmb{x}_2^\top, \cdots, \pmb{x}_T^\top\right]^\top$$

$\pmb{x}_1$
$\pmb{x}_2$
$\vdots$
$\pmb{x}_T$

Latent variables:
$$\pmb{z} = \left[\pmb{z}_1^\top, \pmb{z}_2^\top, \cdots, \pmb{z}_T^\top\right]^\top$$

$\pmb{z}_1$
$\pmb{z}_2$
$\vdots$
$\pmb{z}_T$

- Emission $p.d.f.$ of the observation sequence vector

$$p(\pmb{x}_{1:T} \mid \pmb{z}_{1:T}) = \prod_{t=1}^{T} p(\pmb{x}_t \mid \pmb{z}_t)$$
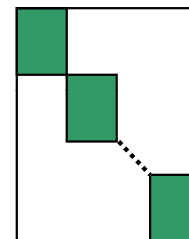
Modeled w/ a single Gaussian distribution

$$= \mathcal{N}\left(\begin{bmatrix} \pmb{x}_1 \\ \pmb{x}_2 \\ \vdots \\ \pmb{x}_T \end{bmatrix}; \begin{bmatrix} \pmb{W} & & & \\ & \pmb{W} & & \\ & & \ddots & \\ & & & \pmb{W} \end{bmatrix} \begin{bmatrix} \pmb{z}_1 \\ \pmb{z}_2 \\ \vdots \\ \pmb{z}_T \end{bmatrix}, \begin{bmatrix} \pmb{\Sigma} & & & \\ & \pmb{\Sigma} & & \\ & & \ddots & \\ & & & \pmb{\Sigma} \end{bmatrix}\right)$$

$$= \mathcal{N}\left(\pmb{x}; \tilde{\pmb{W}}\pmb{z}, \tilde{\pmb{\Sigma}}\right)$$

$$= p(\pmb{x} \mid \pmb{z})$$

$\tilde{\pmb{W}}$

$\tilde{\pmb{\Sigma}}$

App: 1

- Transition $p.d.f.$ of the latent variable sequence vector $z$

$$p(z_{1:T}) = p(z_1)\prod_{t=2}^{T} p(z_t \mid z_{t-1})$$

Mean vector

$$= \mathcal{N}\left(\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_T \end{bmatrix}; \begin{bmatrix} I & & & \\ & A & & \\ & & \ddots & \\ & & & A \end{bmatrix}\begin{bmatrix} \mu_0 \\ z_1 \\ \vdots \\ z_{T-1} \end{bmatrix}, \begin{bmatrix} P_0 & & & \\ & \Gamma & & \\ & & \ddots & \\ & & & \Gamma \end{bmatrix}\right)$$

Subtraction of mean vector from $z$ (*i.e.*, $z - \{$mean vector$\}$) :

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_T \end{bmatrix} - \begin{bmatrix} I & & & \\ & A & & \\ & & \ddots & \\ & & & A \end{bmatrix}\begin{bmatrix} \mu_0 \\ z_1 \\ \vdots \\ z_{T-1} \end{bmatrix} = \begin{bmatrix} I & & & \\ -A & I & & \\ & \ddots & \ddots & \\ & & -A & I \end{bmatrix}\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_T \end{bmatrix} - \begin{bmatrix} \mu_0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Same variables

$$p(\boldsymbol{z}_{1:T}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{I} & & & \\ -\boldsymbol{A} & \boldsymbol{I} & & \\ & \ddots & \ddots & \\ & & -\boldsymbol{A} & \boldsymbol{I} \end{bmatrix}\begin{bmatrix} \boldsymbol{z}_1 \\ \boldsymbol{z}_2 \\ \vdots \\ \boldsymbol{z}_T \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_0 \\ \boldsymbol{0} \\ \vdots \\ \boldsymbol{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{P}_0 & & & \\ & \boldsymbol{\Gamma} & & \\ & & \ddots & \\ & & & \boldsymbol{\Gamma} \end{bmatrix}\right)$$

$$\underbrace{\qquad}_{\tilde{\boldsymbol{A}}} \qquad \underbrace{}_{\tilde{\boldsymbol{\mu}}_0} \qquad \underbrace{\qquad}_{\tilde{\boldsymbol{\Gamma}}}$$

$$= \mathcal{N}\left(\tilde{\boldsymbol{A}}\boldsymbol{z}; \tilde{\boldsymbol{\mu}}_0, \tilde{\boldsymbol{\Gamma}}\right)$$

$$= (2\pi)^{-DT/2}\left|\tilde{\boldsymbol{\Gamma}}\right|^{-1/2}\exp\left(-\frac{1}{2}\left(\tilde{\boldsymbol{A}}\boldsymbol{z} - \tilde{\boldsymbol{\mu}}_0\right)^{\mathsf{T}}\tilde{\boldsymbol{\Gamma}}^{-1}\left(\tilde{\boldsymbol{A}}\boldsymbol{z} - \tilde{\boldsymbol{\mu}}_0\right)\right)$$

$\boxed{\left|\tilde{\boldsymbol{A}}\right| = 1}$

$$= (2\pi)^{-DT/2}\left|\tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{\Gamma}}\tilde{\boldsymbol{A}}^{-\mathsf{T}}\right|^{-1/2}\exp\left(-\frac{1}{2}\left(\boldsymbol{z} - \tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{\mu}}_0\right)^{\mathsf{T}}\tilde{\boldsymbol{A}}^{\mathsf{T}}\tilde{\boldsymbol{\Gamma}}^{-1}\tilde{\boldsymbol{A}}\left(\boldsymbol{z} - \tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{\mu}}_0\right)\right)$$

$$= \mathcal{N}\left(\boldsymbol{z}; \tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{\mu}}_0, \ \tilde{\boldsymbol{A}}^{-1}\tilde{\boldsymbol{\Gamma}}\tilde{\boldsymbol{A}}^{-\mathsf{T}}\right)$$
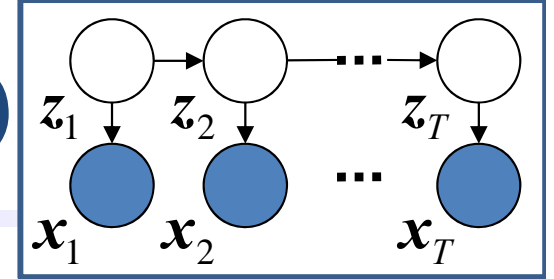
$$= p(\boldsymbol{z})$$

Mean vector

Covariance matrix

# Likelihood Function: $p.d.f.\ p(\boldsymbol{x}_{1:T})$



- $p.d.f.$ of the observation sequence vector $\boldsymbol{x}$:
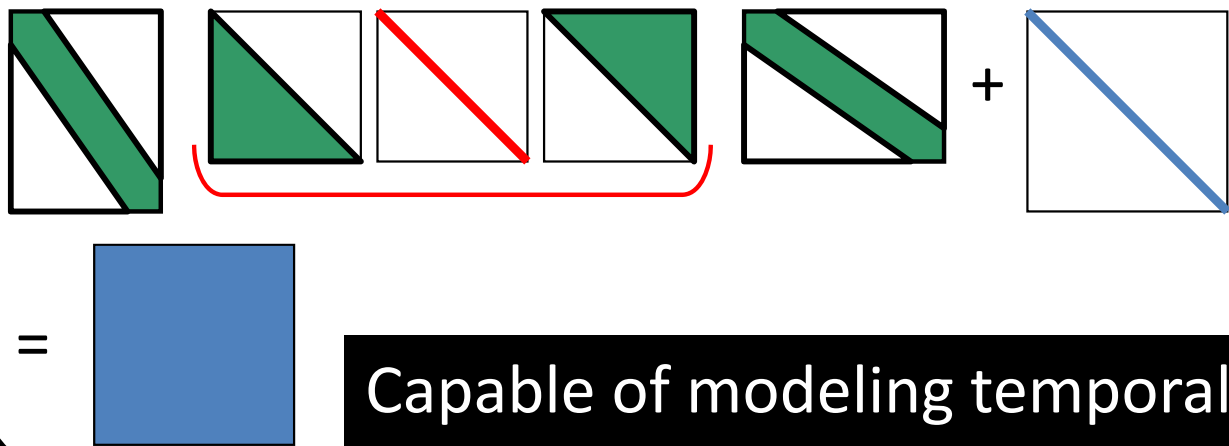
$$p(\boldsymbol{x}) = \int p(\boldsymbol{x} \mid \boldsymbol{z}) p(\boldsymbol{z}) \,\mathrm{d}\boldsymbol{z}$$

$$= \int \mathcal{N}\left(\boldsymbol{x}; \widetilde{\boldsymbol{W}}\boldsymbol{z}, \widetilde{\boldsymbol{\Sigma}}\right) \mathcal{N}\left(\boldsymbol{z}; \underbrace{\widetilde{\boldsymbol{A}}^{-1}\widetilde{\boldsymbol{\mu}}_0}, \underbrace{\widetilde{\boldsymbol{A}}^{-1}\widetilde{\boldsymbol{\Gamma}}\widetilde{\boldsymbol{A}}^{-\top}}\right)\mathrm{d}\boldsymbol{z}$$

$$= \mathcal{N}\left(\boldsymbol{x};\ \underbrace{\widetilde{\boldsymbol{W}}\widetilde{\boldsymbol{A}}^{-1}\widetilde{\boldsymbol{\mu}}_0},\ \underbrace{\widetilde{\boldsymbol{W}}\widetilde{\boldsymbol{A}}^{-1}\widetilde{\boldsymbol{\Gamma}}\widetilde{\boldsymbol{A}}^{-\top}\widetilde{\boldsymbol{W}}^{\top} + \widetilde{\boldsymbol{\Sigma}}}\right)$$

$$\begin{cases} \boldsymbol{x} = \left[\boldsymbol{x}_1^{\top}, \boldsymbol{x}_2^{\top}, \cdots, \boldsymbol{x}_T^{\top}\right]^{\top} \\ \boldsymbol{z} = \left[\boldsymbol{z}_1^{\top}, \boldsymbol{z}_2^{\top}, \cdots, \boldsymbol{z}_T^{\top}\right]^{\top} \end{cases}$$

Mean vector



Covariance matrix



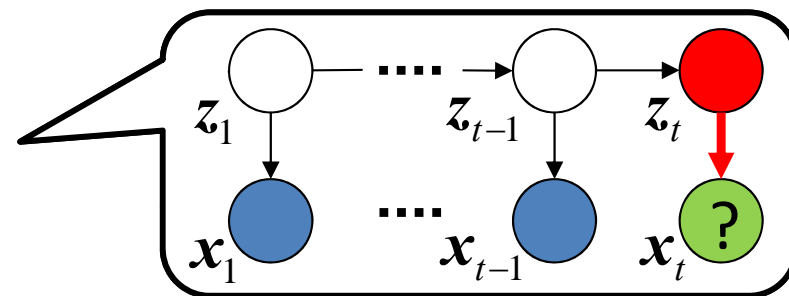Capable of modeling temporal correlation over an sequence

App: 4

# Appendix

# Derivation of $p.d.f.$s in Kalman Filtering
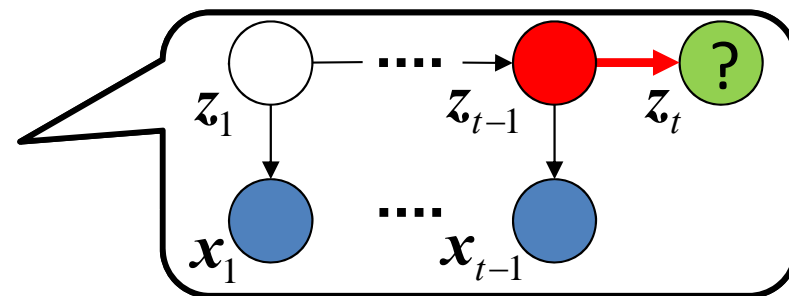
# Forward Algorithm (Kalman Filtering)

- Likelihood function factorized into conditional *p.d.f.*s

$$p(\boldsymbol{x}_{1:T}) = p(\boldsymbol{x}_1)p(\boldsymbol{x}_2 \mid \boldsymbol{x}_1)p(\boldsymbol{x}_3 \mid \boldsymbol{x}_{1:2})p(\boldsymbol{x}_4 \mid \boldsymbol{x}_{1:3})\cdots p(\boldsymbol{x}_t \mid \boldsymbol{x}_{1:t-1})\cdots p(\boldsymbol{x}_T \mid \boldsymbol{x}_{1:T-1})$$

$$p(\boldsymbol{x}_t \mid \boldsymbol{x}_{1:t-1}) = \int p(\boldsymbol{x}_t, \boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1})\,\mathrm{d}\boldsymbol{z}_t$$

$$= \int p(\boldsymbol{x}_t \mid \boldsymbol{z}_t)p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1})\,\mathrm{d}\boldsymbol{z}_t$$

$$p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1}) = \int p(\boldsymbol{z}_t, \boldsymbol{z}_{t-1} \mid \boldsymbol{x}_{1:t-1})\,\mathrm{d}\boldsymbol{z}_{t-1}$$

$$= \int p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1})p(\boldsymbol{z}_{t-1} \mid \boldsymbol{x}_{1:t-1})\,\mathrm{d}\boldsymbol{z}_{t-1}$$

$$p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t}) \propto p(\boldsymbol{x}_t, \boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1})$$

$$\underbrace{\phantom{p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t})}}_{\alpha(\boldsymbol{z}_t)} = p(\boldsymbol{x}_t \mid \boldsymbol{z}_t)p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1})$$

$$= p(\boldsymbol{x}_t \mid \boldsymbol{z}_t)\int p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1})\underbrace{p(\boldsymbol{z}_{t-1} \mid \boldsymbol{x}_{1:t-1})}_{\alpha(\boldsymbol{z}_{t-1})}\,\mathrm{d}\boldsymbol{z}_t$$



App:5

# Derivation of $p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1})$



- **Predicted distribution on the state space**

$$p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t-1}) = \int p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1}) p(\boldsymbol{z}_{t-1} \mid \boldsymbol{x}_{1:t-1}) \, \mathrm{d}\boldsymbol{z}_t$$

Transition $p.d.f. = \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{A}\boldsymbol{z}_{t-1}, \boldsymbol{\Gamma})$
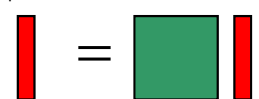
Assumed to be $\mathcal{N}(\boldsymbol{z}_{t-1}; \boldsymbol{\mu}_{t-1}, \boldsymbol{P}_{t-1})$
(Its derivation will be given later.)

$$= \int \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{A}\boldsymbol{z}_{t-1}, \boldsymbol{\Gamma}) \mathcal{N}(\boldsymbol{z}_{t-1}; \boldsymbol{\mu}_{t-1}, \boldsymbol{P}_{t-1}) \, \mathrm{d}\boldsymbol{z}_{t-1}$$

$$= \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{A}\boldsymbol{\mu}_{t-1}, \boldsymbol{A}\boldsymbol{P}_{t-1}\boldsymbol{A}^{\mathsf{T}} + \boldsymbol{\Gamma})$$

$$= \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{\mu}_{t|t-1}, \boldsymbol{P}_{t|t-1})$$

Predicted mean $\quad: \quad \boldsymbol{\mu}_{t|t-1} = \boldsymbol{A}\boldsymbol{\mu}_{t-1}$



Predicted covariance $: \quad \boldsymbol{P}_{t|t-1} = \boldsymbol{A}\boldsymbol{P}_{t-1}\boldsymbol{A}^{\mathsf{T}} + \boldsymbol{\Gamma}$



$$\boldsymbol{\mu}_{1|0} = \boldsymbol{\mu}_0$$
$$\boldsymbol{P}_{1|0} = \boldsymbol{P}_0$$

App:6

# Derivation of $p(\boldsymbol{x}_t | \boldsymbol{x}_{1:t-1})$



- **Likelihood function**

  (*i.e.*, predicted distribution on the observation space)

$$p(\boldsymbol{x}_t | \boldsymbol{x}_{1:t-1}) = \int p(\boldsymbol{x}_t | \boldsymbol{z}_t) p(\boldsymbol{z}_t | \boldsymbol{x}_{1:t-1}) \, d\boldsymbol{z}_t$$

Emission $p.d.f. = \mathcal{N}(\boldsymbol{x}_t; \boldsymbol{W}\boldsymbol{z}_t, \boldsymbol{\Sigma})$

Predicted distribution $= \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{\mu}_{t|t-1}, \boldsymbol{P}_{t|t-1})$

$$= \int \mathcal{N}(\boldsymbol{x}_t; \boldsymbol{W}\boldsymbol{z}_t, \boldsymbol{\Sigma}) \mathcal{N}(\boldsymbol{z}_t; \boldsymbol{\mu}_{t|t-1}, \boldsymbol{P}_{t|t-1}) \, d\boldsymbol{z}_t$$

$$= \mathcal{N}(\boldsymbol{x}_t; \boldsymbol{W}\boldsymbol{\mu}_{t|t-1}, \boldsymbol{W}\boldsymbol{P}_{t|t-1}\boldsymbol{W}^\top + \boldsymbol{\Sigma})$$

Mean : $\boldsymbol{W}\boldsymbol{\mu}_{t|t-1}$

Covariance: $\boldsymbol{W}\boldsymbol{P}_{t|t-1}\boldsymbol{W}^\top + \boldsymbol{\Sigma}$

# Derivation of $\alpha(z_t) = p(z_t \mid x_{1:t})$



- **Updated distribution on the state space**

$$\alpha(z_t) \propto p(x_t \mid z_t) p(z_t \mid x_{1:t-1})$$

> Posterior $\propto$ Likelihood x Prior

$$= p(x_t \mid z_t) \int \left\{ p(z_t \mid z_{t-1}) \alpha(z_{t-1}) \right\} dz_{t-1}$$

$$\begin{cases} p(z_t \mid z_{t-1}) = \mathcal{N}(z_t; A z_{t-1}, \Gamma) \\ p(x_t \mid z_t) = \mathcal{N}(x_t; W z_t, \Sigma) \end{cases}$$

Assuming that $\alpha(z_t) = \mathcal{N}(z_t; \mu_t, P_t)$

$$\mathcal{N}(z_t; \mu_t, P_t) \propto \mathcal{N}(x_t; W z_t, \Sigma) \int \left\{ \underbrace{\mathcal{N}(z_t; A z_{t-1}, \Gamma) \mathcal{N}(z_{t-1}; \mu_{t-1}, P_{t-1})}_{\mathcal{N}(z_{t-1}; ? z_t + ?, ?) \mathcal{N}(z_t; ?, ?)} \right\} dz_{t-1}$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\mathcal{N}(z_t; ?, ?)}$$

Kalman gain matrix : $K_t = P_{t|t-1} W^\top \left( W P_{t|t-1} W^\top + \Sigma \right)^{-1}$

Updated mean : $\mu_t = \mu_{t|t-1} + K_t \left( x_t - W \mu_{t|t-1} \right)$

> Error between predicted and observed data

Updated covariance : $P_t = (I - K_t W) P_{t|t-1}$

# Appendix

# Derivation of $p.d.f.$s in Kalman Smoothing

# Backward Algorithm (Kalman Smoothing)

- Joint posterior *p.d.f.* of $\boldsymbol{z}_t$ and $\boldsymbol{z}_{t+1}$



$$p(\boldsymbol{z}_t, \boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T}) = p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T})\,p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t+1}, \boldsymbol{x}_{1:T})$$

$$= p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T})\,p(\boldsymbol{z}_t \mid \boldsymbol{z}_{t+1}, \boldsymbol{x}_{1:t})$$

$$= p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T})\,\frac{p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t})\,p(\boldsymbol{z}_{t+1} \mid \boldsymbol{z}_t, \boldsymbol{x}_{1:t})}{p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:t})}$$

$$= p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t})\,\frac{p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T})\,p(\boldsymbol{z}_{t+1} \mid \boldsymbol{z}_t)}{p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:t})}$$

Transition *p.d.f.*

Updated *p.d.f.* in forward algorithm

Predicted *p.d.f.* in forward algorithm

- Posterior *p.d.f.* of $\boldsymbol{z}_t$, *i.e.*, $p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:T}) = \int p(\boldsymbol{z}_t, \boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T})\,\mathrm{d}\boldsymbol{z}_{t+1}$

$$\underbrace{p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:T})}_{\gamma(\boldsymbol{z}_t)} = \underbrace{p(\boldsymbol{z}_t \mid \boldsymbol{x}_{1:t})}_{\alpha(\boldsymbol{z}_t)} \int \frac{\overbrace{p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:T})}^{\gamma(\boldsymbol{z}_{t+1})}\,p(\boldsymbol{z}_{t+1} \mid \boldsymbol{z}_t)}{p(\boldsymbol{z}_{t+1} \mid \boldsymbol{x}_{1:t})}\,\mathrm{d}\boldsymbol{z}_{t+1}$$

# Derivation of $\gamma(z_t) = p(z_t \,|\, x_{1:T})$

$$\gamma(z_t) = \alpha(z_t) \int \frac{\gamma(z_{t+1}) p(z_{t+1} \,|\, z_t)}{p(z_{t+1} \,|\, x_{1:t})} \, dz_{t+1}$$

Assuming that $\gamma(z_t) = \mathcal{N}\left(z_t; \hat{\boldsymbol{\mu}}_t, \hat{\boldsymbol{P}}_t\right)$

$$\mathcal{N}\left(z_t; \hat{\boldsymbol{\mu}}_t, \hat{\boldsymbol{P}}_t\right) = \mathcal{N}\left(z_t; \boldsymbol{\mu}_t, \boldsymbol{P}_t\right) \int \underbrace{\frac{\mathcal{N}\left(z_{t+1}; \hat{\boldsymbol{\mu}}_{t+1}, \hat{\boldsymbol{P}}_{t+1}\right) \mathcal{N}\left(z_{t+1}; \boldsymbol{A} z_t, \boldsymbol{\Gamma}\right)}{\mathcal{N}\left(z_{t+1}; \boldsymbol{\mu}_{t+1|t}, \boldsymbol{P}_{t+1|t}\right)}}_{\underbrace{\mathcal{N}\left(z_{t+1}; ? z_t + ?, ?\right) \mathcal{N}\left(z_t; ?, ?\right)}_{\mathcal{N}\left(z_t; ?, ?\right)}} \, dz_{t+1}$$

$$\boldsymbol{J}_t = \boldsymbol{P}_t \boldsymbol{A}^{\mathsf{T}} \boldsymbol{P}_{t+1|t}^{-1}$$

Smoothed mean : $\hat{\boldsymbol{\mu}}_t = \boldsymbol{\mu}_t + \boldsymbol{J}_t \left(\hat{\boldsymbol{\mu}}_{t+1} - \boldsymbol{\mu}_{t+1|t}\right)$

Smoothed covariance : $\hat{\boldsymbol{P}}_t = \boldsymbol{P}_t + \boldsymbol{J}_t \left(\hat{\boldsymbol{P}}_{t+1} - \boldsymbol{P}_{t+1|t}\right) \boldsymbol{J}_t^{\mathsf{T}}$

# Tips: Matrix Inversion Lemma

- Condition 1: Matrix $A$ and its inverse matrix $A^{-1}$ are given.

$$A = \boxed{\phantom{XX}} \qquad A^{-1} = \boxed{\phantom{XX}}^{-1}$$

- Condition 2: Fluctuation generated on a lower-dimensional subspace $vNv^{\mathrm{T}}$ is added to the matrix $A$.

$$vNv^{\mathrm{T}} = \qquad A + vNv^{\mathrm{T}} = \quad + $$

- Under these conditions, an inverse matrix $\left(A + vNv^{\mathrm{T}}\right)^{-1}$ can be calculated as follows:

$$\left(A + vNv^{\mathrm{T}}\right)^{-1} = A^{-1} - A^{-1}v\left(N^{-1} + v^{\mathrm{T}}A^{-1}v\right)^{-1}v^{\mathrm{T}}A^{-1}$$

Calculation of an inverse matrix on a lower dimensional space