

Methods of Learning Dynamical Systems

Apr.19, 2018

Takehisa YAIRI (矢入健久)

E-mail: yairi@ailab.t.u-tokyo.ac.jp

What is this lecture about ?

- Machine learning for system identification
 - In machine learning, it is called "learning dynamical systems"
- Why "learning dynamical" systems ?
 - Methods for learning static systems (e.g., supervised classification and regression) are already matured
 - Most of existing systems are dynamic in nature
 - Inference methods for dynamical systems are also matured
- Classified into three approaches:
 1. Maximum likelihood (EM-based) approach
 2. Spectral (subspace) approach
 3. (Deep) neural network approach

Schedule

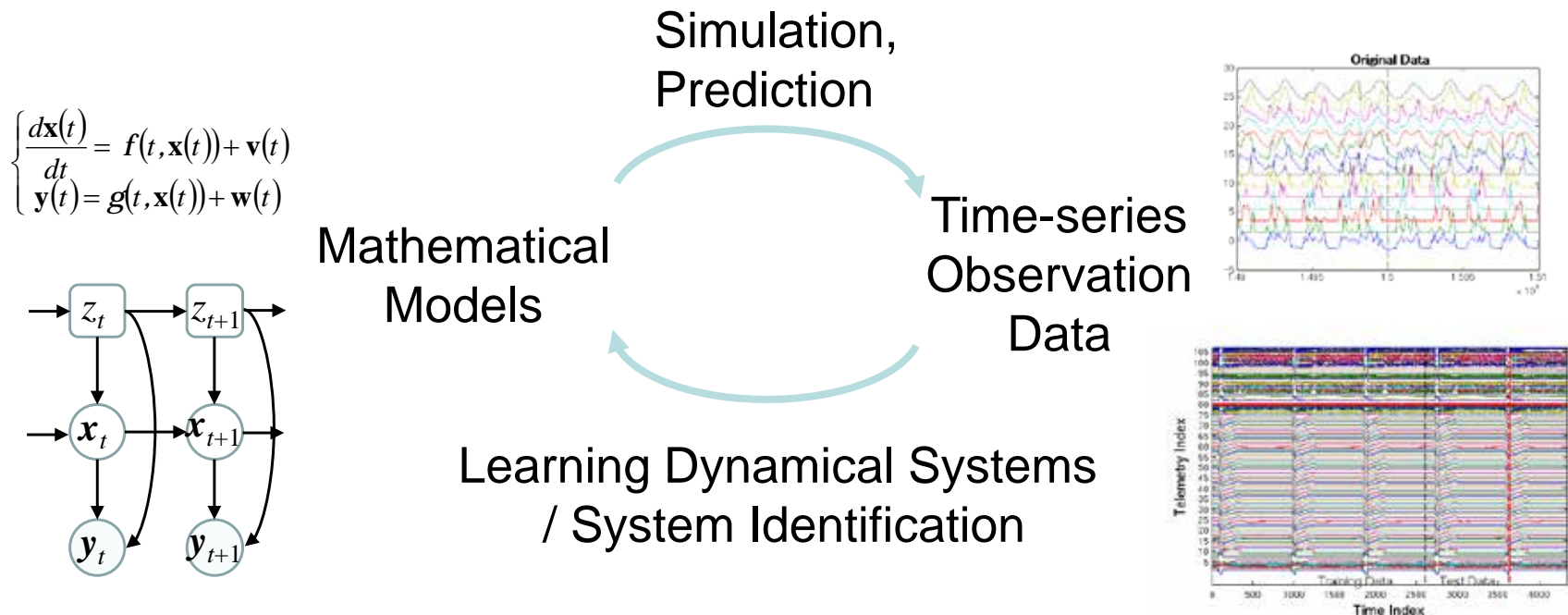
- 4/19: Guidance & Introduction :
 - What is "learning dynamical systems" ?
- 5/10: Maximum Likelihood Approach 1 :
 - EM algorithm for linear dynamical systems
- 5/24: Maximum Likelihood Approach 2 :
 - Learning switching linear systems
- 5/30: Spectral Approach 1 :
 - Subspace identification
- 6/21: Spectral Approach 2 :
 - Non-linearization by kernel, Mixture model
- 7/5: Neural Network Approach :
 - Deep learning for dynamical systems

Prof. Hori will give lectures on 4/5, 4/26, 5/17, 6/14, 6/28, 7/12

Introduction to Learning Dynamical Systems

Learning Dynamical Systems

- Estimating the models of unknown dynamical systems from time-series observation data
- Known as "System Identification" in the control theory



State Space Model (SSM) in Machine Learning

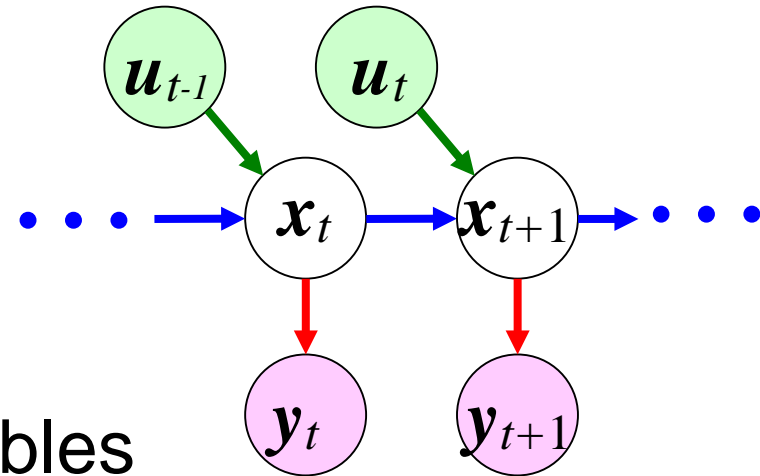
- SSM is popular in machine learning, as well as in control
- Regarded as a special case of latent variable models (LVM)
 - SSM \approx LVM with temporal structure (dynamics)
- Probabilistic representation is often used
 - Can be unified with hidden Markov models (HMM)

State	$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}_t$	\iff	$p(\mathbf{x}_{t+1} \mathbf{x}_t, \mathbf{u}_t)$
transition	$\mathbf{w}_t \sim N(\mathbf{0}, \mathbf{Q})$		$= N(\mathbf{x}_{t+1} f(\mathbf{x}_t, \mathbf{u}_t), \mathbf{Q})$

Observation	$\mathbf{y}_t = h(\mathbf{x}_t) + \mathbf{v}_t$	\iff	$p(\mathbf{y}_t \mathbf{x}_t)$
(Output)	$\mathbf{v}_t \sim N(\mathbf{0}, \mathbf{R})$		$= N(\mathbf{y}_t h(\mathbf{x}_t), \mathbf{R})$

State Space Model (SSM) in Machine Learning (Cont.)

- Graph representation
 - Bayesian Networks [Pearl 88]
 - Nodes : random variables
 - Directed edges : dependencies
- Continuous and discrete variables can be mixed
 - HMM : discrete latent variables
 - Linear dynamical systems (LDS) : continuous latent variables
 - Switching LDS, Dynamic Bayesian Networks : both continuous and discrete state variables



Methods of Learning Dynamical Systems

We classifies existing methods into three categories:

1. Maximum likelihood estimation approach

- Iteration of state estimation and model estimation

2. Spectral approach

- Inspired by subspace identification

3. Neural network approach

- A recent trend

Note that they are not necessarily mutually exclusive

1. Maximum Likelihood Approach (In General)

- Not limited to MLE in the narrow sense
- Intended to include maximum a posteriori (MAP) and Bayesian estimation

- MLE : $\hat{\theta}_{ML} = \arg \max_{\theta} p(\mathbf{y}_{1:T} | \theta) = \arg \max_{\theta} \int p(\mathbf{y}_{1:T}, \mathbf{x}_{1:T} | \theta) d\mathbf{x}_{1:T}$

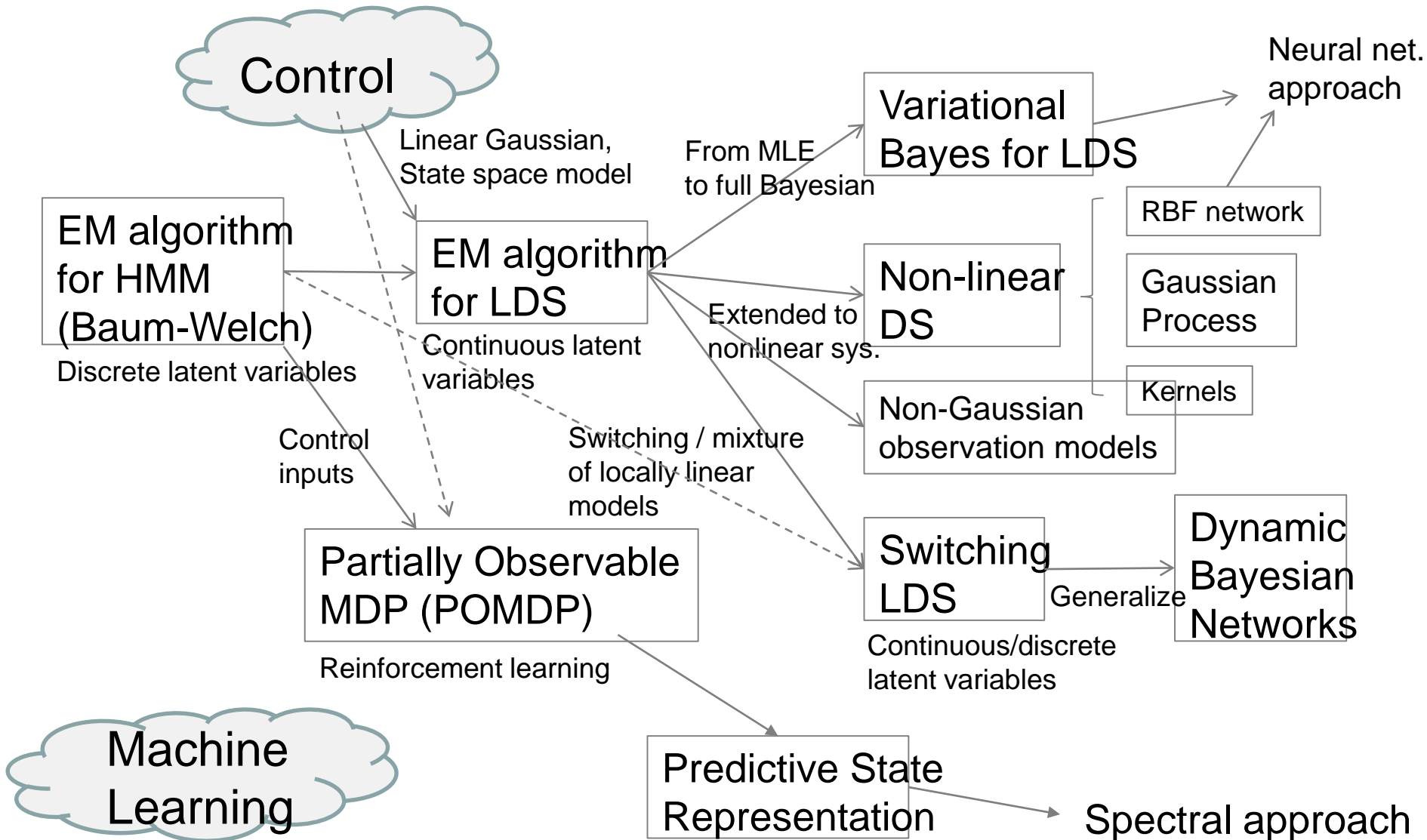
- MAP: $\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | \mathbf{y}_{1:T}) = \arg \max_{\theta} p(\mathbf{y}_{1:T} | \theta) p(\theta)$

- Bayesian: $\hat{\alpha} = \arg \max_{\alpha} p(\mathbf{y}_{1:T} | \alpha)$
 $= \arg \max_{\alpha} \left\{ \int p(\mathbf{y}_{1:T} | \theta) p(\theta | \alpha) d\theta \right\}$

$\mathbf{y}_{1:T}$: Observations
 $\mathbf{x}_{1:T}$: States
 θ : Parameters
 α : Hyper param

- EM algorithm is the base method
 - In Bayesian estimation, variational inference and MCMC are also employed
- Iterations of state estimation and model estimation

Pedigree of Maximum Likelihood Approach



Expectation Maximization Algorithm

- A mandatory technique for machine learning researchers
- Maximum likelihood (or maximum a posteriori) estimation method for latent variable models
- More efficient than gradient methods
- Initialization is critical due to local minima

Given:

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

Repeat until convergence

1. [E-step] Compute posterior dist. $q^*(X) = p(X | Y, \Theta^{(t)})$
and $Q(\Theta | \Theta^{(t)}) = E_{q^*(X)}[\ln p(Y, X | \Theta)]$
2. [M-step] Maximize $Q(\Theta | \Theta^{(t)})$ w.r.t. Θ
 $\Theta^{(t+1)} \leftarrow \arg \max_{\Theta} Q(\Theta | \Theta^{(t)})$
Expected complete log-likelihood
3. $t \leftarrow t + 1$

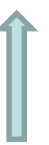
We want to maximize
 $p(Y | \Theta) = \int p(Y, X | \Theta) dX$
w.r.t. Θ but, X is also unknown



Compute $q(X) = p(X | Y, \Theta)$



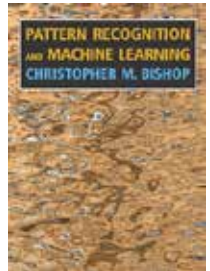
Iteration



Maximize
 $E_{q(X)}[\ln p(Y, X | \Theta)]$
w.r.t. Θ

Learning Hidden Markov Model (HMM) by EM Algorithm

- Starting point of learning dynamical systems (maybe..)
 - Known as Baum-Welch algorithm
 - Included in popular machine learning textbooks such as PRML [Bishop 06]
 - HMM is popular in natural language processing, bioinformatics, activity recognition, etc. (but not in control)
- E-step : Fix model parameters, then compute the posterior distribution of discrete latent variables
- M-step : Fix the posterior of latent variables, then maximize expected likelihood w.r.t. parameters



Learning Linear Dynamical Systems by EM algorithm

- Fundamental idea is almost the same with HMM
 - If the model is known, states can be estimated
 - If the states are known, model can be estimated
- Firstly introduced by [Ghahramani & Hinton 96]
 - Technical report, not a journal nor conference paper
- Re-introduced in PRML [Bishop 06]
- E-step = Rauch-Tung-Striebel (RTS) smoothing
 - In machine learning, it is known as Kalman smoothing
 - In fact, "Kalman filtering" is often referred as a general inference technique for state space models

1. MLE approach

EM Algorithm for Linear Dynamical Systems (Summary)

Model: $\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{w}_t$ $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$ (*) Control input \mathbf{u}_t is not considered for simplicity

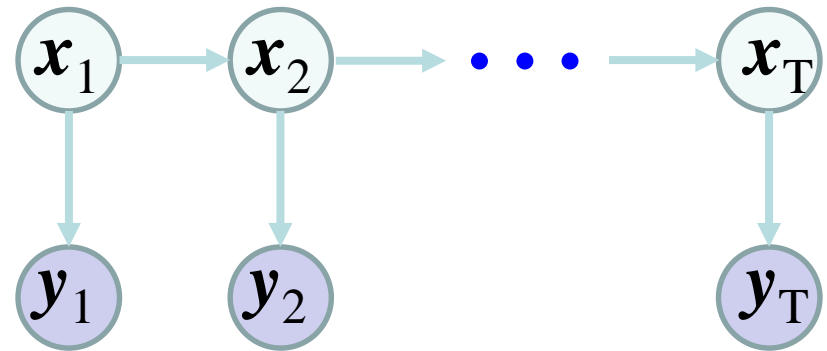
$\mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{v}_t$ $\mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$

Given :

- Obs. sequence : $\mathbf{y}_{1:T}$

Find :

- System matrices : \mathbf{A}, \mathbf{C}
- Noise covariance: \mathbf{Q}, \mathbf{R}
- Posterior dist. of state sequence $\mathbf{x}_{1:T}$: $p(\mathbf{x}_{1:T} | \mathbf{y}_{1:T})$
 - Assume Gaussians : $p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{V}_t)$
- Prior dist. of initial state: $p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_0, \mathbf{V}_0)$



EM Algorithm for Linear Dynamical Systems (Summary)

- Initialize estimates of model parameters

- Repeat until convergence:

[E-step] Fix model parameters, then compute posterior of state sequence

For $t=1:T-1$ Kalman filter

$$p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t+1}) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{m}_{t+1}, \mathbf{V}_{t+1}) \quad \left\{ \begin{array}{l} \mathbf{P}_t = \mathbf{A}\mathbf{V}_t\mathbf{A}^T + \mathbf{Q} \quad \mathbf{K}_{t+1} = \mathbf{P}_t\mathbf{C}^T(\mathbf{C}\mathbf{P}_t\mathbf{C}^T + \mathbf{R})^{-1} \\ \mathbf{m}_{t+1} = \mathbf{A}\mathbf{m}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{C}\mathbf{A}\mathbf{m}_t) \\ \mathbf{V}_{t+1} = (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{C})\mathbf{P}_t \end{array} \right.$$

For $t=T-1:1$ RTS smoother

$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_t | \hat{\mathbf{m}}_t, \hat{\mathbf{V}}_t) \quad \left\{ \begin{array}{l} \mathbf{J}_t = \mathbf{V}_t\mathbf{A}^T\mathbf{P}_t^{-1} \\ \hat{\mathbf{m}}_t = \mathbf{m}_t + \mathbf{J}_t(\hat{\mathbf{m}}_{t+1} - \mathbf{A}\mathbf{m}_t) \\ \hat{\mathbf{V}}_t = \mathbf{J}_t\hat{\mathbf{V}}_{t+1}\mathbf{J}_t^T + \mathbf{V}_t - \mathbf{J}_t\mathbf{A}\mathbf{V} = \mathbf{V}_t + \mathbf{J}_t(\hat{\mathbf{V}}_{t+1} - \mathbf{P}_t)\mathbf{J}_t^T \end{array} \right.$$

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

[M-step] Fix posterior of state sequence, then update the model parameters

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

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

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

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

From MLE to MAP and Bayesian Estimation

- Maximum a posteriori estimation (MAPE):
Set prior distribution on the model parameters A, B, Q, R , then maximize the posterior probability w.r.t. the parameters
 - Reasonable when we want to use some prior knowledge or when data is not sufficient
 - EM algorithm can be used (if the prior distribution is conjugate)
- Bayesian estimation : Find the posterior distribution of the parameters $p(A, B, Q, R | y_{1:T})$
 - Uncertainty of estimation is taken into consideration
 - Can be used for model selection based on marginal likelihood
 - Solved by variational Bayes[Barber06] and MCMC

Non-Gaussian Observation Model

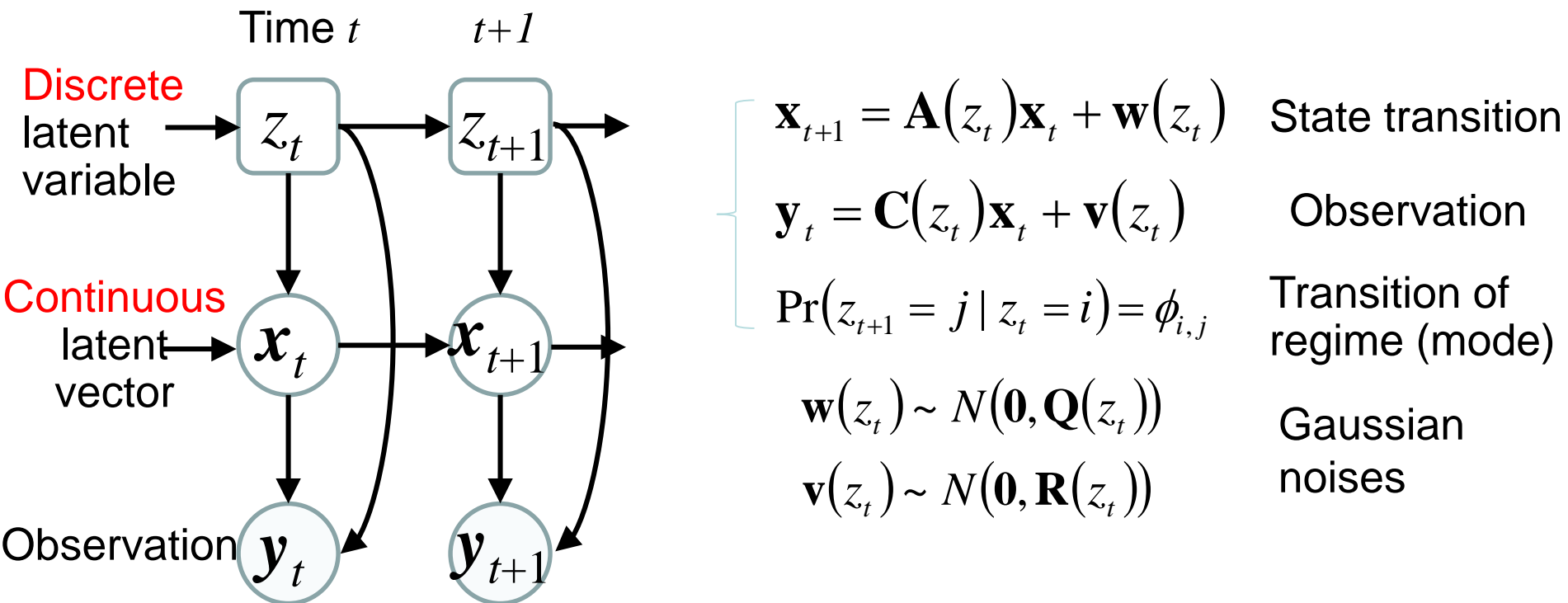
- Advantage of maximum likelihood approach :
Can be applied to "any" state space model
 - Not limited to linear Gaussian models
 - But, E-step and M-step can be more complicated..
- For example,
 - [Macke 11] [Gao 15][Park 15] :
Dynamics of neuron activity are learned from data.
Linear Gaussian state transition model + Poisson
distribution observation model. EM algorithm (or
variational EM) is used.

1. MLE approach

Switching Linear Dynamical System (SLDS)

[Murphy 98][Ghahramani&Hinton 00]

Switching among several linear models stochastically



- A hybrid of HMM and linear dynamical system
- E-step is approximately computed, as analytical solution is intractable

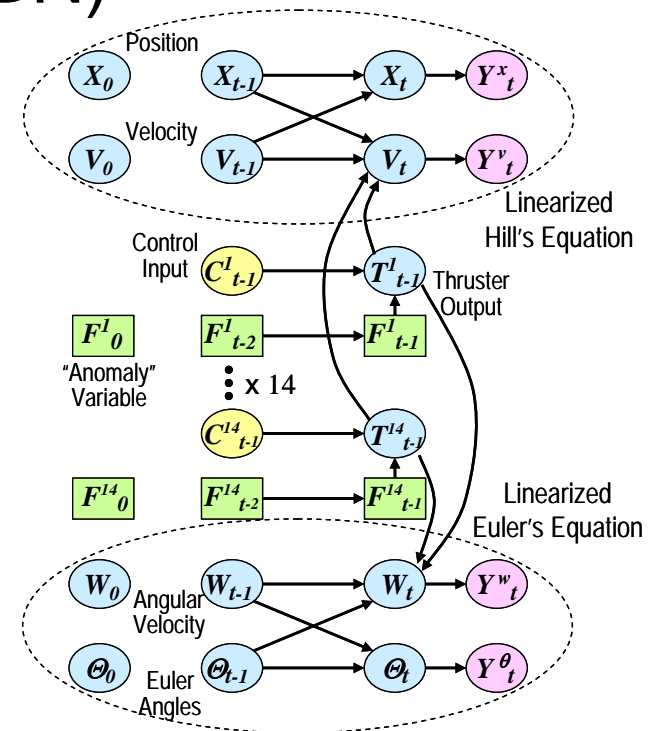
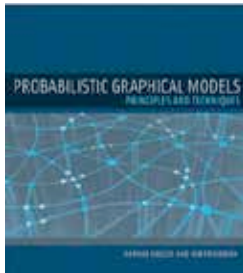
Learning of Switching Linear Dynamical Systems

- Still actively studied
- [Fox 2008] “Nonparametric Bayesian Learning of Switching Linear Dynamical Systems”, NIPS-2008
 - Dirichlet process clustering is applied
 - All parameters are estimated by MCMC
- [Chiappa 2008] “Using Bayesian Dynamical Systems for Motion Template Libraries”, NIPS-2008
 - Parameters are estimated by variational Bayes
- Applied to many purposes such as activity recognition

1. MLE approach

Dynamic Bayesian Networks

- HMM, LDS, SLDS and more complicated state space models are generalized into Dynamic Bayesian Networks (DBN)
- DBN is a special case of Bayesian networks
- Text book of DBN : [Koller 09]



Non-linear State Space Model with Radial Basis Function Network

- [Ghahramani & Roweis 1999] “Learning Nonlinear Dynamical Systems using the EM Algorithm”
- Non-linear state model f and observation model g are approximated by radial basis function networks

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}_t \quad \text{State transition model}$$

$$\approx \sum_{i=1}^I h_i \rho_i(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \mathbf{w}_t$$

Gaussian RBF

- State sequence and models are iteratively estimated by EM algorithm
 - E-step : Extended Kalman (RTS) smoothing

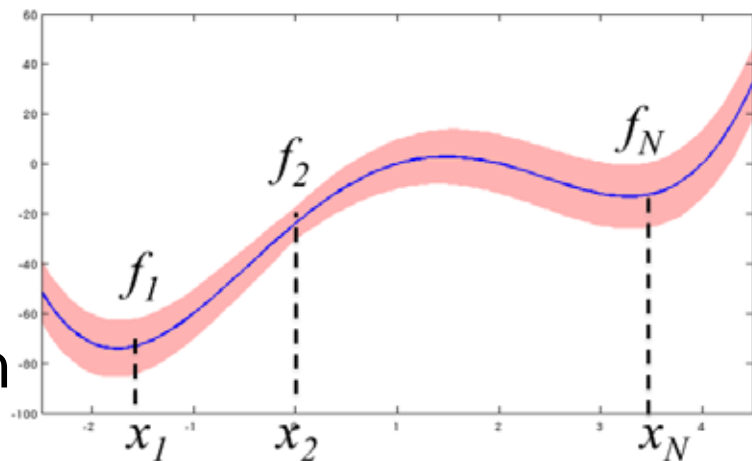
Non-linear State Space Model with Kernels

Kernel Kalman Filter (KKF) [Ralaivola 05]

- Nonlinearization by kernels (RKHS)
 - Assumed that linear dynamics and linear observation model can hold in the feature space
 - Kernel PCA is used to obtain bases in the feature space
- Learned (estimated) by EM algorithm
- Preimage is necessary
- Similar ideas of using RKHS can be seen even recently [Zhu 14]

Non-linear State Space Model with Gaussian Process

- Gaussian process regression [Rasmussen 06]
 - Supervised learning of $p(y | x)$
 - Gaussian process :
Prior distribution over a function
 $y = f(x)$
 - If $\{y_1, y_2, \dots, y_N, y_t\}$ is Gaussian, then
 $p(y_t | y_1, \dots, y_N)$ is also Gaussian
 - More compatible with probabilistic models (than SVM(SVR))
 - "Sparse" GP techniques are also available
- Gaussian Process Latent Variable Model [Lawrence 03]
 - Unsupervised learning (nonlinear dimensionality reduction)



$$p(y) = \int p(y | x) p(x) dx$$

y : High dimensional observation

x : Low dimensional latent vector

Gaussian Process State Space Models(1)

- Supervised learning of SSM
 - Assumption : State sequence $X=\{x_1, \dots, x_N\}$ is directly observable !
 - Reduced to supervised regression problem
 - State and observation models are learned by ordinary Gaussian process regression
- GP-Bayesfilter[Ko 08], GP-ADF[Deisenroth 09]
 - $f(x_t, u_t)$ and $g(x_t)$ are learned by GP regression
 - Semi-parametric
 - Sparse GP is also considered
 - Applied to system identification of blimp dynamics

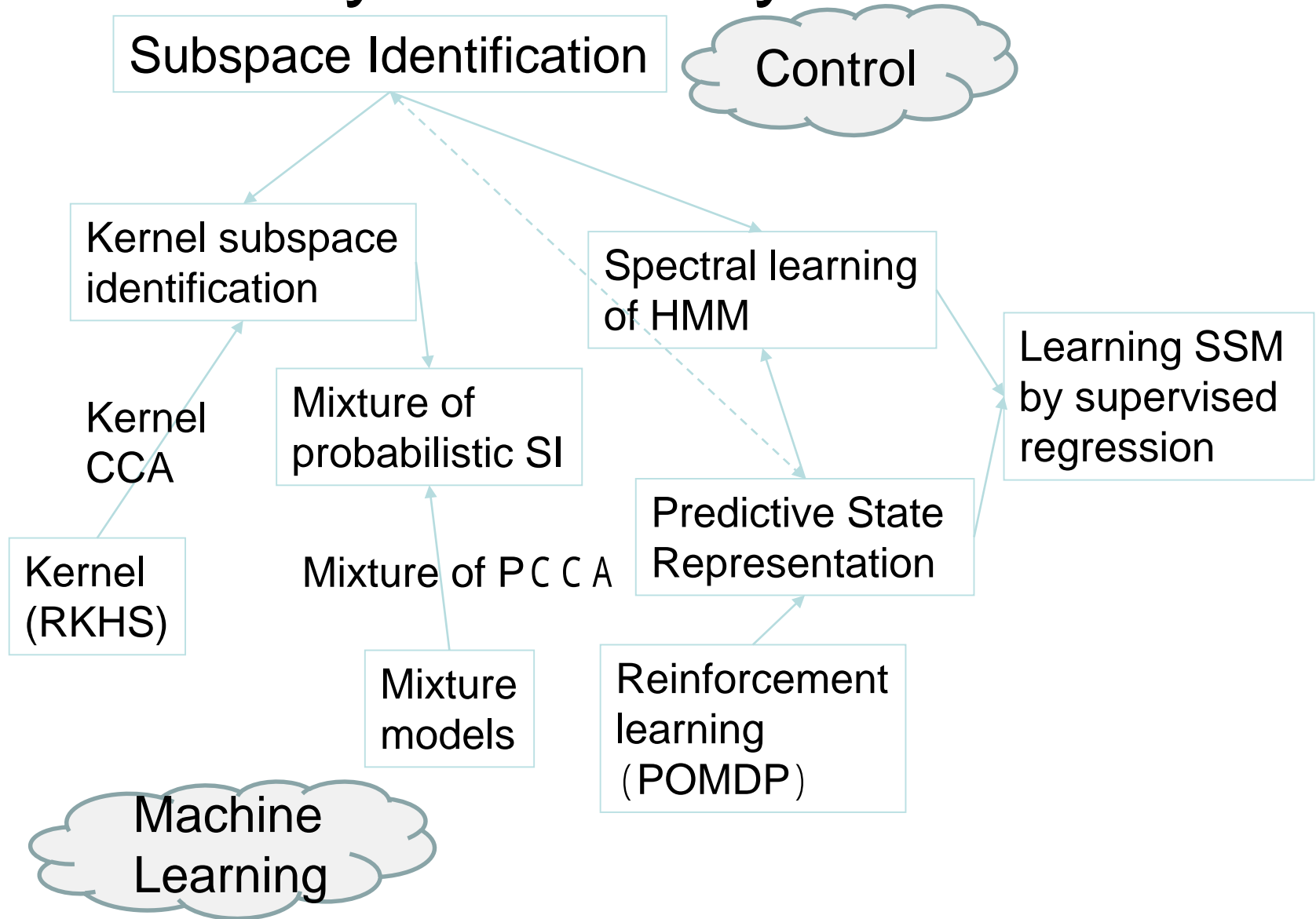
Gaussian Process State Space Models(2)

- Gaussian Process Dynamical Models [Wang 08]
 - $f(\mathbf{x}_t, \mathbf{u}_t)$ and $g(\mathbf{x}_t)$ are learned by GP regression
 - Marginal likelihood is maximized w.r.t. state sequence $\mathbf{x}_{1:T}$ and hyper parameters
 - Initialized by (ordinary) PCA
- GP Inference and Learning [Turner 10]
 - Sparse GP model using pseudo inputs
- Variational GPDS [Damianou 2011]
 - Variational Bayes inference instead of EM algorithm
- Variational GPSSM [Frigola 2013]
 - Combination of GP-based state model and parametric observation model

2. Spectral Learning of Dynamical Systems (Overview)

- Subspace identification from a machine learning perspective
- Why is it called "spectral" learning ?
 - Solved by eigen-decomposition and SVD
 - Related to manifold learning (?)
- Non-iterative algorithm and global optimum
- Less flexible than the maximum likelihood approach
- Spectral learning for HMM [Hsu 09]

Pedigree of Spectral Learning of Dynamical Systems



2. Spectral approach

Kernel Subspace Identification [Kawahara 06]

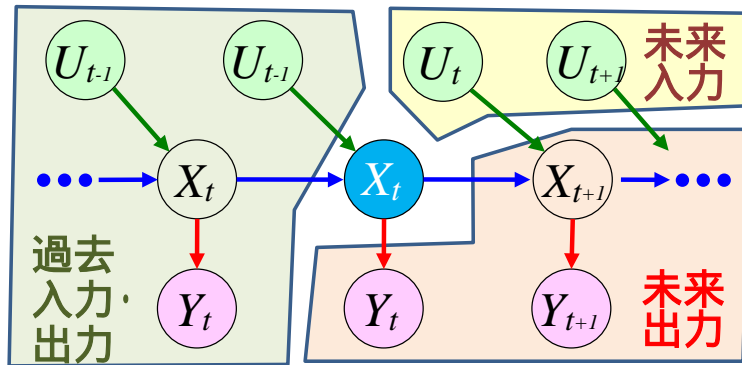
Y. Kawahara, T. Yairi, and K. Machida, “A kernel subspace method by stochastic realization for learning nonlinear dynamical systems”, NIPS-2006

- Non-linearization of subspace identification based on canonical correlation analysis [Larimore 90][Katayama 99] by kernel (RKHS)
 - Kernel canonical correlation analysis [Akaho 01][Bach & Jordan 02]
 - KKF[Ralaivola 05] : learned by EM algorithm
- Pre-image problem is inevitable
- Pioneer work of introducing subspace identification to machine learning community

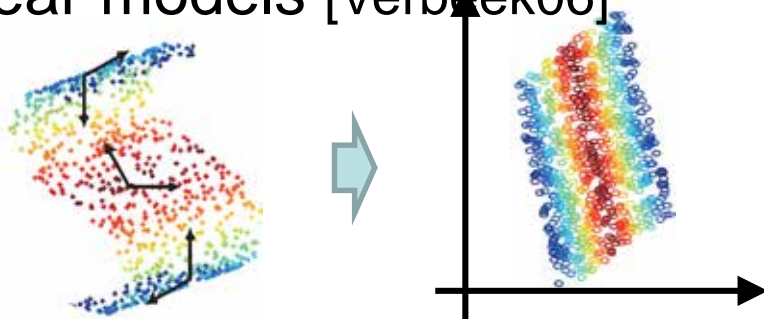
2. Spectral approach

Mixture of Probabilistic Subspace Identification [Joko 11]

Subspace Identification by CCA
[Larimore 90][Katayama 99]



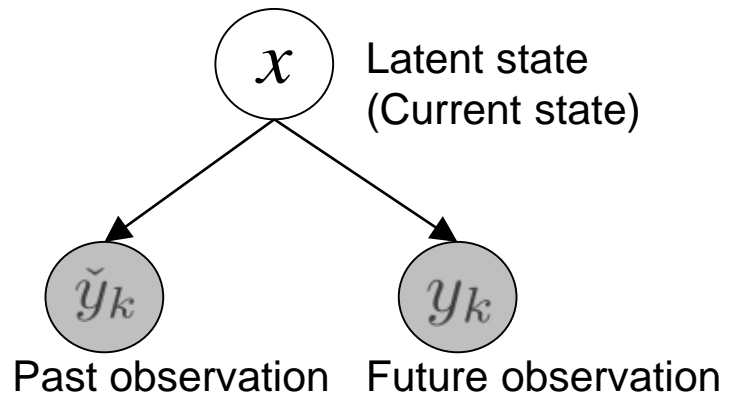
Alignment of locally
linear models [Verboek06]



Multiple local coordinates

Global coordinates

Probabilistic interpretation
of CCA [Bach&Jordan 06]



Mixture

Non-linearization of CCA
subspace identification by
mixture of Probabilistic CCA

- Similar to SLDS model

2. Spectral approach

Spectral Learning of HMM [Hsu 09]

Daniel J. Hsu, Sham M. Kakade, and Tong Zhang, “A Spectral Algorithm for Learning Hidden Markov Models”, COLT 2009.

- For a long time, EM (Baum-Welch) algorithm was believed to be the only way to learn HMM
- Inspired by subspace identification
 - Consider the canonical correlation between past and future observation
 - (Latent) state sequence and transition/output probabilities are implicitly computed
- Limitations
 - Limited to discrete observations
 - Assumption of one-step observability

Spectral Learning of HMM (Cont.)

The seminal work of [Hsu 09] was rapidly extended

- Frustration to EM algorithm
- [Siddiqi 10] Continuous observations
- [Song 10] Non-Gaussian continuous observation
- [Anandkumar 12] Generalization as a "method of moments"
- [Subakan 14] Mixture of HMM
- [Zhang 15] Latent states with tree-like structure
- [Kandasamy 16] Non-parametric observation model

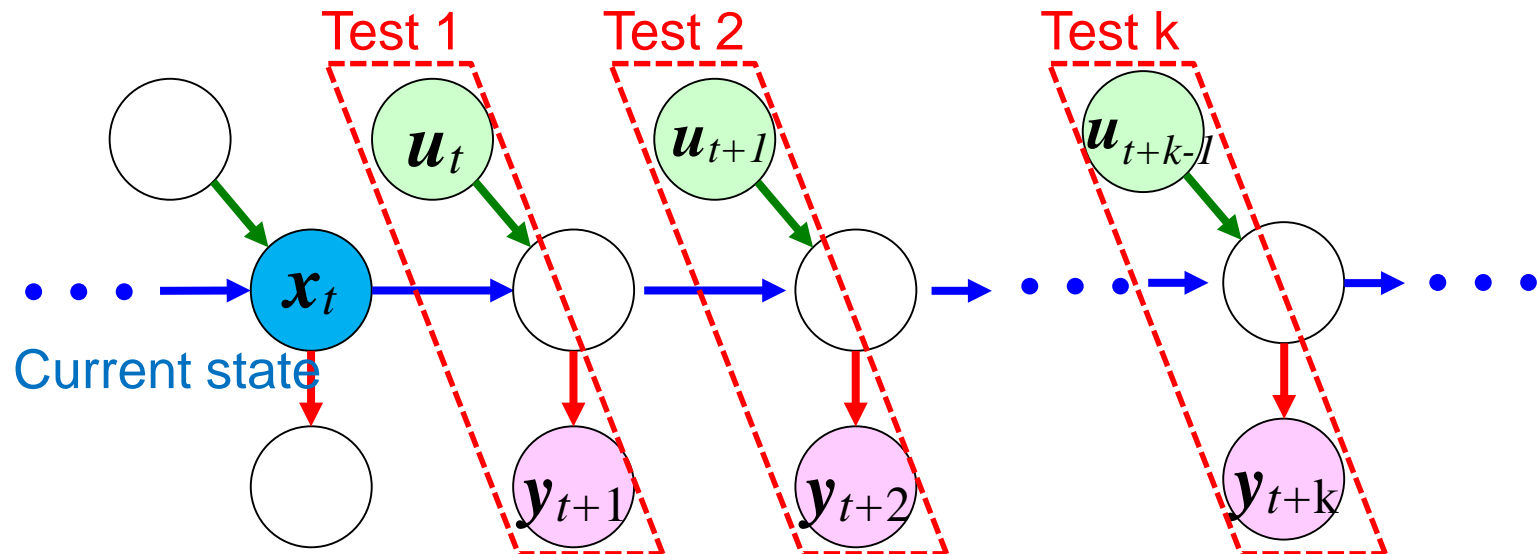
Spectral Learning of Continuous State Space Models

- Spectral learning of HMM is "re-imported" to continuous state space models
 - [Buesing 12] Extension of Ho-Kalman realization algorithm to Poisson observation model
- Another idea : Use the past observation sequence, instead of estimating latent state explicitly
 - Reduced to supervised regression problems
 - [Langford 09], [Hefny 2015], [Sun 2016]
 - Predictive State Representation (PSR) [Littman 01]

2. Spectral approach

Predictive State Representation (PSR)

- Originally developed as a state representation for partially observable environment [Littman 01] [Singh 04]
 - Extension of Observable operator models (OOM) [Jaeger 00]
- Instead of estimating the current state, predict a set of tests (pairs of input and output) in future
 - Predicting test results in future \approx Guessing the current state



2. Spectral approach

Predictive State Representation (Cont.)

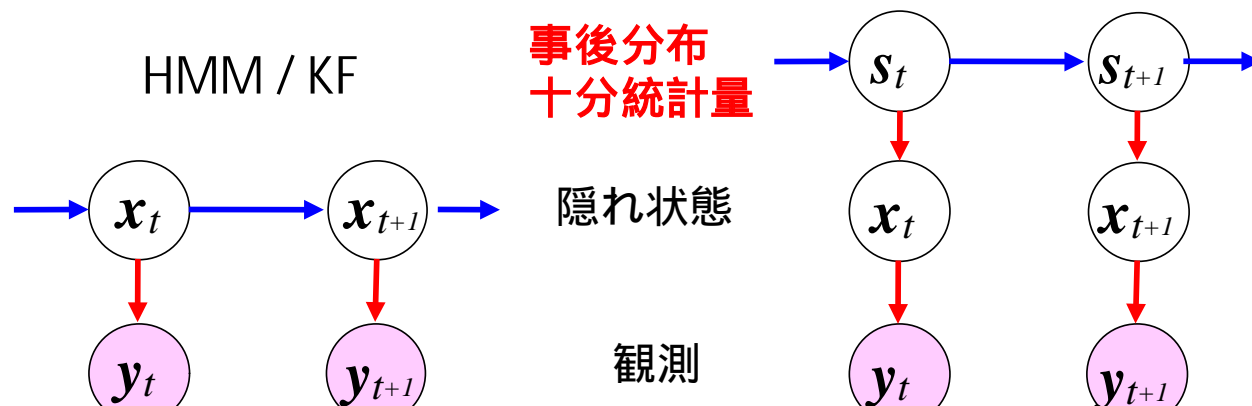
- Sufficient statistic for future test results \approx (current) state
- Predicting future test results based on past results \approx State estimation (filtering)
- Transformed PSR [Rosencrantz 04] : Obtain a minimum set of bases necessary to predict any future test results
- Close relation to canonical variate [Akaike 75], subspace identification[Boots 09]

2. Spectral approach

Langford's Method

"Learning Nonlinear Dynamic Models“, Langford ,
Salakhutdinov, Zhang, ICML-2009

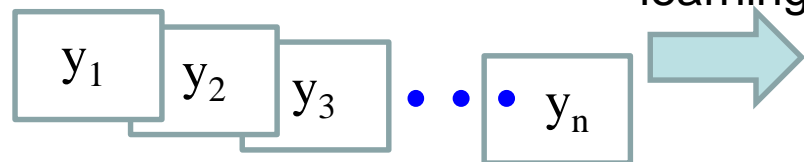
- Convert the learning of stochastic dynamical systems into deterministic supervised regression problems
- Transforms the probabilistic state \mathbf{x}_t into a deterministic variable s_t
 - Sufficient Posterior Representation: SPR_{SPR}



3. Neural Network Approach to Learning Dynamical Systems

- Maybe, [Roweis & Ghahramani 98] is the first
 - State space model with RBFN
- Recurrent neural networks (RNN)
- Deep learning for dynamical systems ?
- Variational autoencoder [Kingma 14]
 - Unsupervised learning of generative latent variable models

High-dimensional data
(e.g., images)



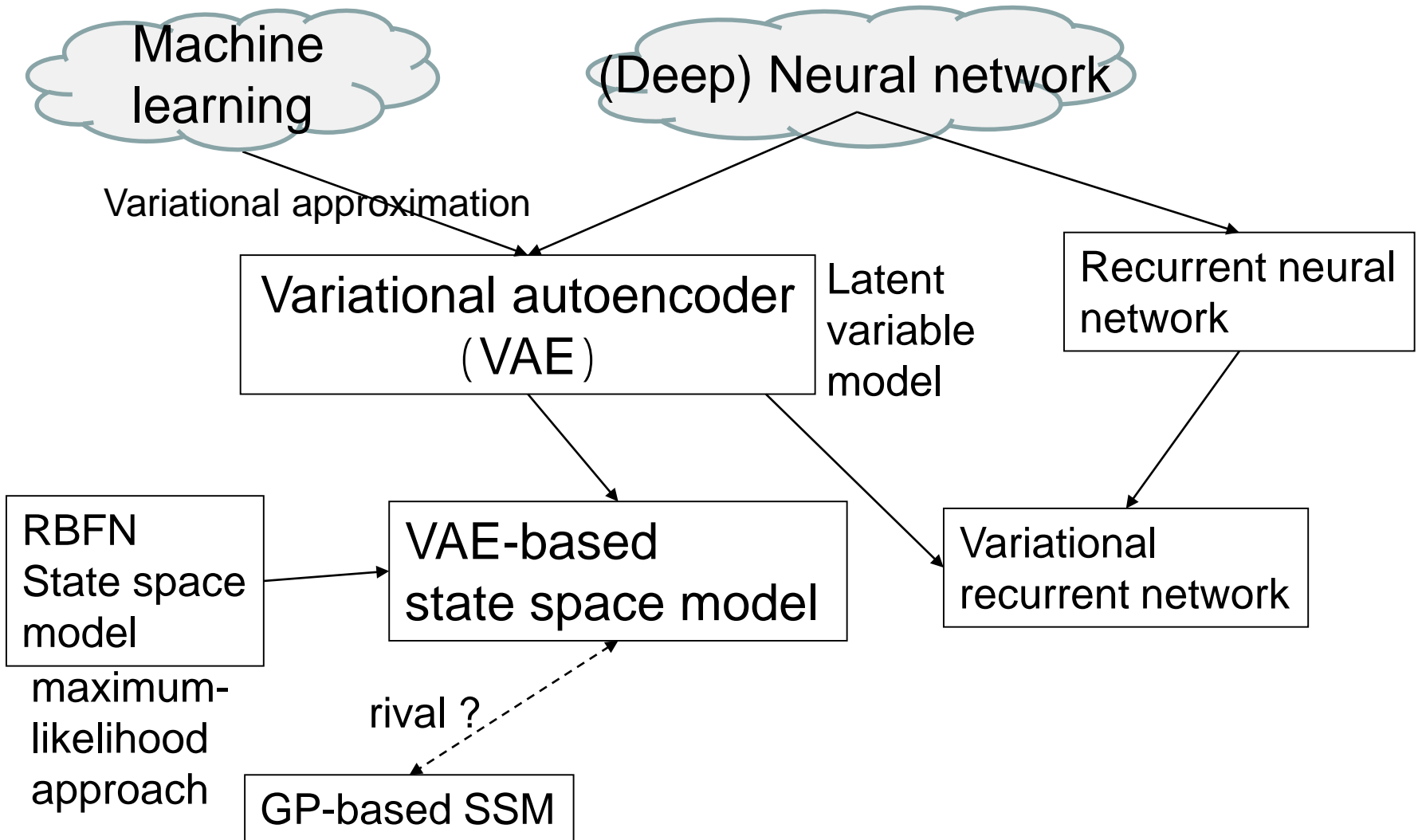
Latent variable model (e.g., VAE)

$$p(\mathbf{y}) = \int p(\mathbf{y} | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

\mathbf{y} : high-dimensional observation

\mathbf{x} : low-dimensional latent state

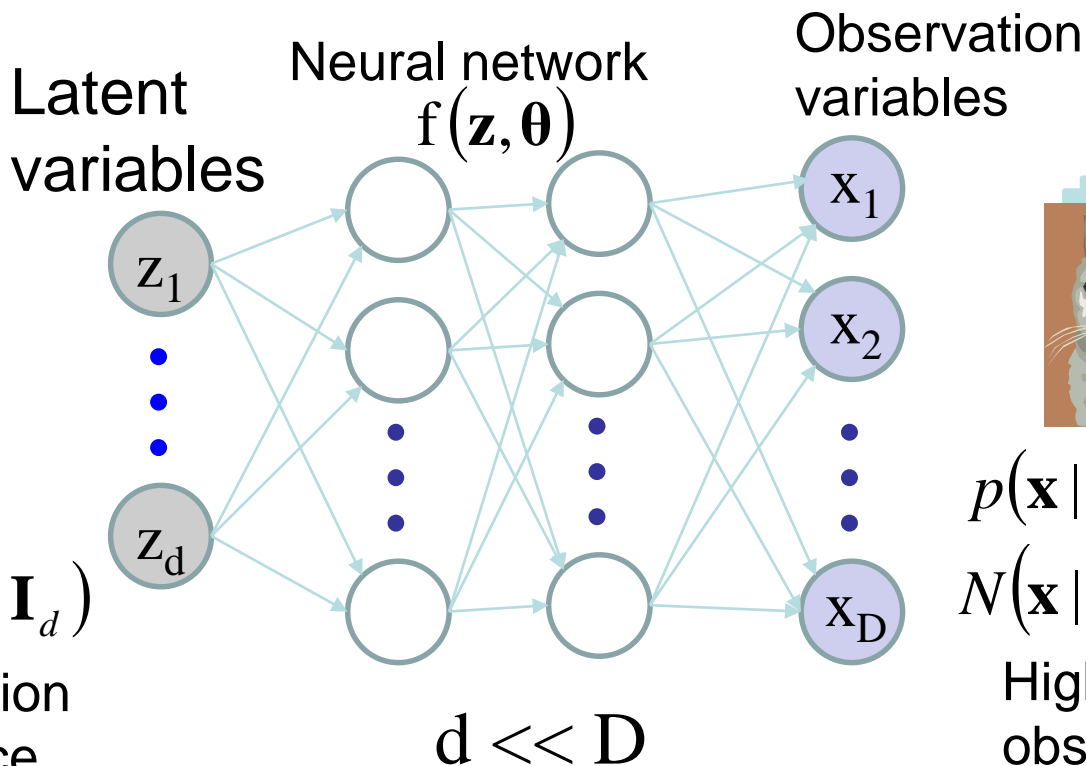
Pedigree of Neural Network Approach to Learning Dynamical Systems



Variational Autoencoder [Kingma 14]

Variational Autoencoder (VAE) [Kingma 14]

- A latent variable model using neural network
 - Probabilistic, generative model
 - Powerful approximation ability of deep neural network



$$p(\mathbf{x} | \mathbf{z}) = N(\mathbf{x} | f(\mathbf{z}, \theta), \sigma^2 \cdot \mathbf{I}_D)$$

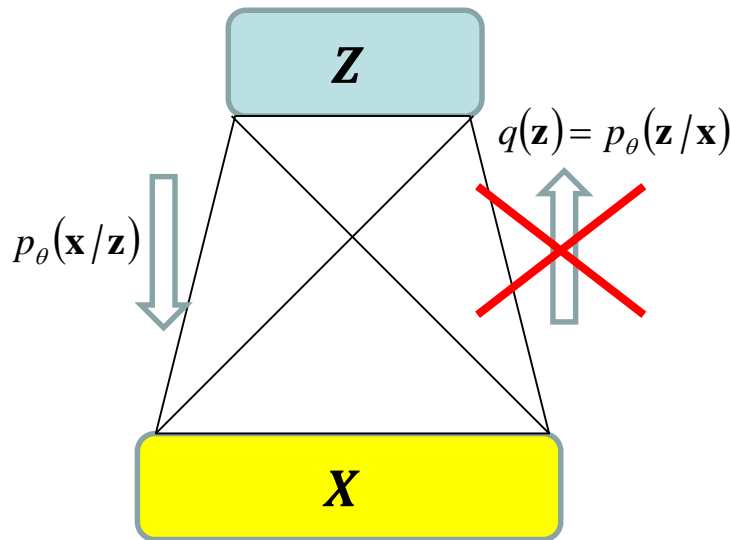
High-dimensional observation data

$$p(\mathbf{z}) = N(\mathbf{z} | \mathbf{0}, \mathbf{I}_d)$$

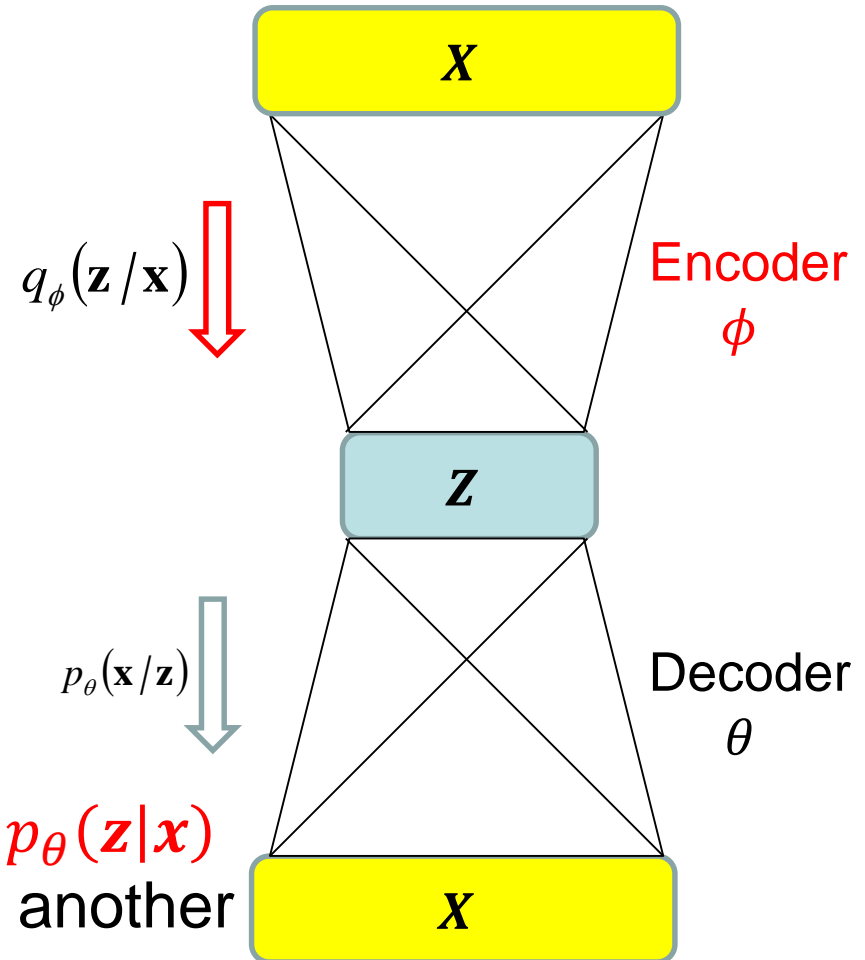
Normal distribution
in low-dim. space

Variational Autoencoder (cont.)

EM algorithm



VAE (Encoder-Decoder)



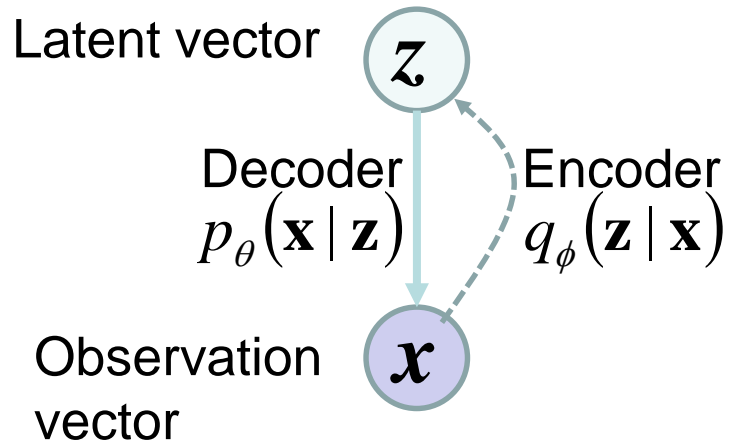
Replace the intractable posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ with its approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$ by another neural network

3. Neural network approach

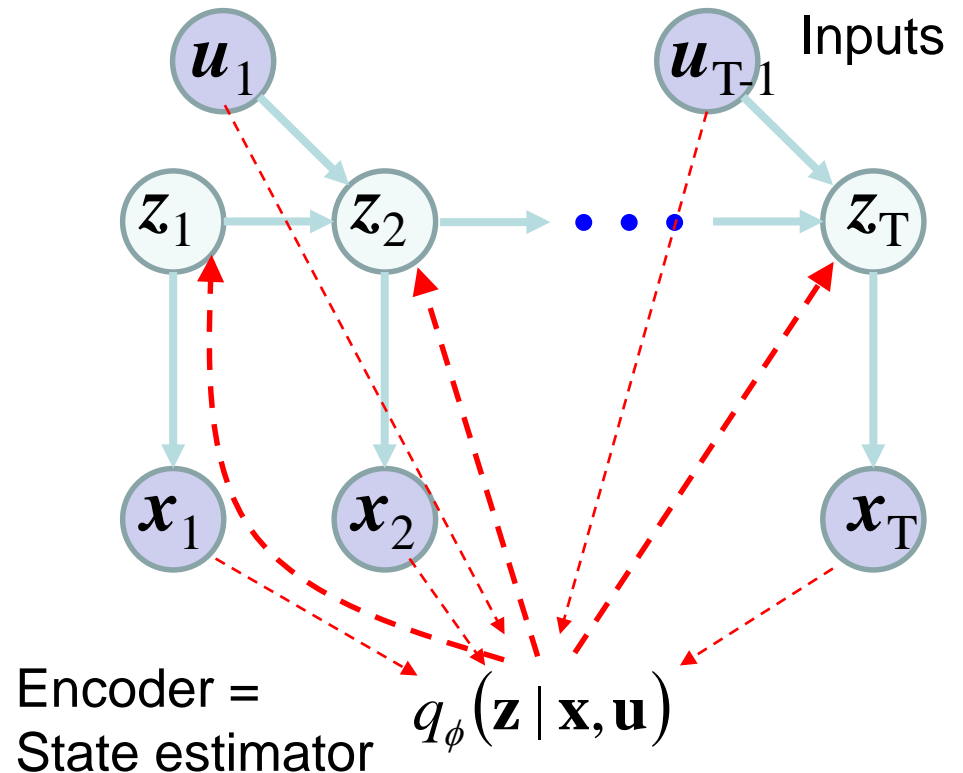
Deep Kalman Filter [Krishnan 15]

State space model with VAE

Variational Auto-Encoder



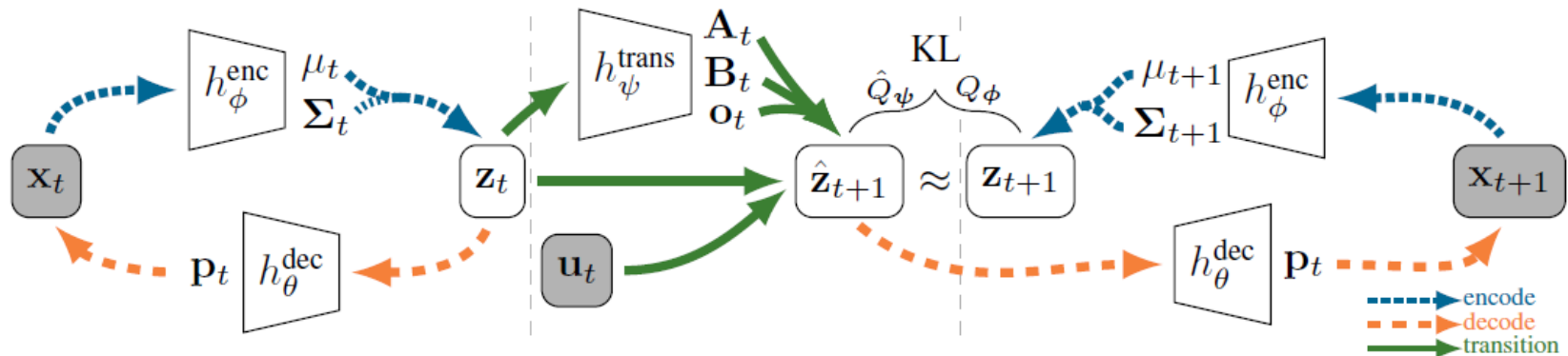
Deep Kalman Filter



3. Neural network approach

Embed to control [Watter 2015]

- Learning a non-linear dynamical system from a sequence of raw images
- Observation function (from latent state to observation) is modeled by VAE
- State transition is assumed to be locally linear

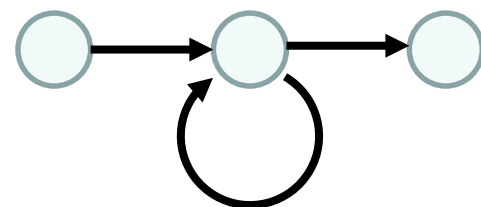


[Watter 2015]

3. Neural network approach

Recurrent Neural Networks

- A neural network model for sequence or time-series data
- RNN has the internal state (memory)
 - Not probabilistic, but deterministic



Inspired by Variational Autoencoder

- Recurrent Latent Variable Model[Chung 15]
 - An alternative to the state space model
- Stochastic recurrent network (STORN)[Bayer 15]

Next Week

- We will begin by maximum likelihood approach
- We will derive the EM algorithm for linear dynamical systems