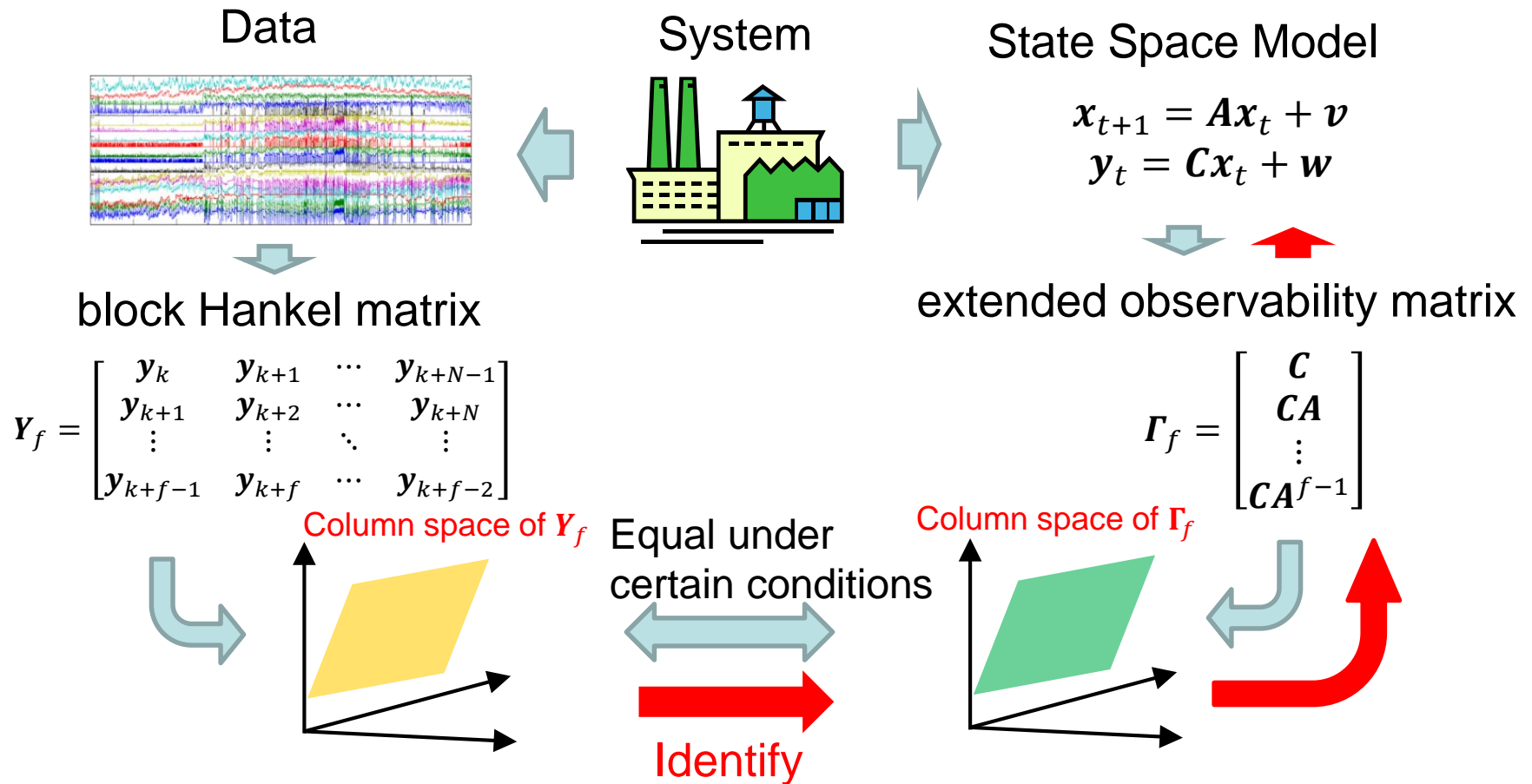# Canonical Correlation Analysis and Spectral Learning
## - *What is the state ?* -

Jul.5, 2018

Takehisa YAIRI （矢入健久）

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

# Review of Subspace Identification

Data

System

State Space Model

$$x_{t+1} = Ax_t + v$$
$$y_t = Cx_t + w$$

block Hankel matrix

extended observability matrix

$$Y_f = \begin{bmatrix} y_k & y_{k+1} & \cdots & y_{k+N-1} \\ y_{k+1} & y_{k+2} & \cdots & y_{k+N} \\ \vdots & \vdots & \ddots & \vdots \\ y_{k+f-1} & y_{k+f} & \cdots & y_{k+f-2} \end{bmatrix}$$

$$\Gamma_f = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{f-1} \end{bmatrix}$$

Column space of $Y_f$

Equal under certain conditions

Column space of $\Gamma_f$

Identify

- ## Advantage:
  - Global optimum is obtained without iteration by linear algebraic operations

2

# Subspace Identification : CCA Approach*

$$Y_f \Pi_{U_f}^\perp = \Gamma_f L_z Z_p \Pi_{U_f}^\perp + G_f E_f$$

Known · Unknown · Known · Residual

Variable set 1 · Variable set 2

Extended state space representation

$$Y_f = \Gamma_f X_k + H_f U_f + G_f E_f$$

States estimated from past inputs and outputs

$$X_k \equiv L_z Z_p$$

Obtained by Canonical Correlation Analysis (CCA):

Define $W_r = \left(Y_f \Pi_{U_f}^\perp Y_f{}^T\right)^{-1/2}$ and $W_c = \left(Z_p \Pi_{U_f}^\perp Z_p{}^T\right)^{-1/2}$

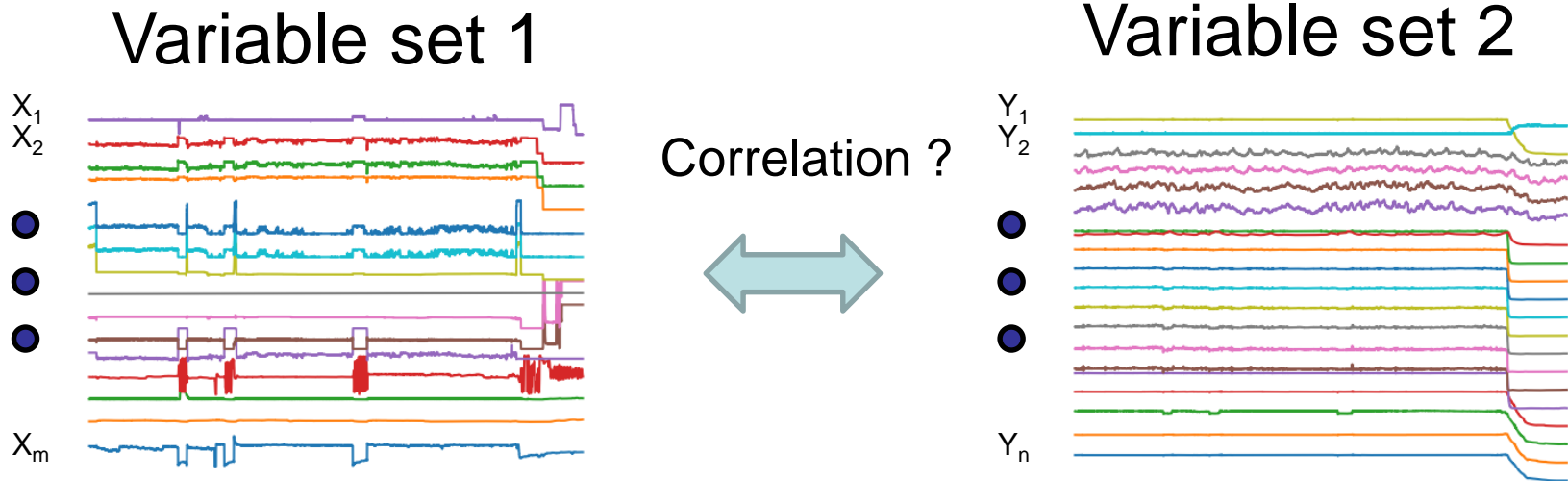Perform SVD on $W_r Y_f \Pi_{U_f}^\perp Z_p{}^T W_c \approx U_d S_d V_d^T$

Then, $\hat{\Gamma}_f \leftarrow W_r{}^{-1} U_d S_d{}^{1/2}$

(*) Also known as canonical variate analysis (CVA)

# Canonical Correlation Analysis (CCA)

# What is CCA ?

Assume there are two sets of variables (signals)

Variable set 1          Variable set 2



Correlation ?

Q : How are the two sets of variables correlated ?

- It's easy to examine the correlation of any pair of variables
- But variables in each set have correlations themselves
- Something like "normalization" and "orthogonalization" is necessary

# Canonical Correlation Analysis (1)

There are two random vectors $x$ and $y$

$$x \in R^m , \quad y \in R^n$$

For simplicity, assume both $x$ and $y$ are "centered", i.e., their means are zero.

$$E[x] = \mathbf{0}_m , \quad E[y] = \mathbf{0}_n$$

Covariance matrices *of* and *between* $x$ and $y$

$$\text{var}(x) = E[xx^T] = \Sigma_{xx} , \text{var}(y) = E[yy^T] = \Sigma_{yy}$$

$$\text{cov}(x, y) = E[xy^T] = \Sigma_{xy} \quad \Rightarrow \quad \Sigma_{yx} = \Sigma_{xy}{}^T$$

# Canonical Correlation Analysis (2)

Think of constructing synthetic variables $u$ and $v$ by linear combinations of $\boldsymbol{x}$ and $\boldsymbol{y}$, respectively.

$$u = \boldsymbol{a}^T\boldsymbol{x} = a_1 x_1 + a_2 x_2 + \cdots + a_m x_m$$

$$v = \boldsymbol{b}^T\boldsymbol{y} = b_1\boldsymbol{y} + b_2 y_2 + \cdots + b_n y_n$$

Problem: Find $\boldsymbol{a}$ and $\boldsymbol{b}$ so that the correlation between $u$ and $v$ is maximized

$$\rho = \mathrm{cor}(u, v) = \frac{\mathrm{cov}(u, v)}{\sqrt{\mathrm{var}(u) \cdot \mathrm{var}(v)}}$$

$$= \frac{\boldsymbol{a}^T\boldsymbol{\Sigma}_{xy}\boldsymbol{b}}{\sqrt{\boldsymbol{a}^T\boldsymbol{\Sigma}_{xx}\boldsymbol{a}}\sqrt{\boldsymbol{b}^T\boldsymbol{\Sigma}_{yy}\boldsymbol{b}}}$$

# Canonical Correlation Analysis (3)

Impose constraints on $\boldsymbol{a}$ and $\boldsymbol{b}$, so that the variances of $u$ and $v$ become 1

$$\text{var}(u) = \text{var}(\boldsymbol{a}^T \boldsymbol{x}) = \boldsymbol{a}^T \boldsymbol{\Sigma}_{xx} \boldsymbol{a} = 1$$

$$\text{var}(v) = \text{var}(\boldsymbol{b}^T \boldsymbol{y}) = \boldsymbol{b}^T \boldsymbol{\Sigma}_{yy} \boldsymbol{b} = 1$$

The problem is formulated as,

$$(\boldsymbol{a}_1, \boldsymbol{b}_1) = \underset{\boldsymbol{a}^T \boldsymbol{\Sigma}_{xx} \boldsymbol{a}=1 \,,\, \boldsymbol{b}^T \boldsymbol{\Sigma}_{yy} \boldsymbol{b}=1}{\text{argmax}} \boldsymbol{a}^T \boldsymbol{\Sigma}_{xy} \boldsymbol{b}$$

It can be solved by Lagrange multiplier, but there is a more elegant method.

# Canonical Correlation Analysis (4)

Let square root matrices of $\Sigma_{xx}$ and $\Sigma_{yy}$ be ${\Sigma_{xx}}^{1/2}$ and ${\Sigma_{yy}}^{1/2}$ , respectively. Then, define $\boldsymbol{c}$ and $\boldsymbol{d}$ as,

$$\boldsymbol{c} = {\boldsymbol{\Sigma}_{xx}}^{1/2}\boldsymbol{a} \quad , \text{ and } \boldsymbol{d} = {\boldsymbol{\Sigma}_{yy}}^{1/2}\boldsymbol{b}$$

Note that $\Sigma_{xx}$ and $\Sigma_{yy}$ are positive definite matrices

The problem turns to be

$$(\boldsymbol{c}_1, \boldsymbol{d}_1) = \underset{\|\boldsymbol{c}\|=1\,,\|\boldsymbol{d}\|=1}{\operatorname{argmax}} \; \boldsymbol{c}^T \left( \boldsymbol{\Sigma}_{xx}^{-T/2} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{xx}^{-1/2} \right) \boldsymbol{d}$$

This problem can be solved by SVD !

$$\boldsymbol{\Sigma}_{xx}^{-T/2} \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{xx}^{-1/2} = \boldsymbol{CSD}^T \quad \Rightarrow \quad \begin{cases} \boldsymbol{c}_1 \leftarrow \text{1st column of } \boldsymbol{C} \\ \boldsymbol{d}_1 \leftarrow \text{1st column of } \boldsymbol{D} \end{cases}$$

$$\Rightarrow \quad \boldsymbol{a}_1 = {\boldsymbol{\Sigma}_{xx}}^{1/2}\boldsymbol{c}_1, \text{ and } \boldsymbol{b}_1 = {\boldsymbol{\Sigma}_{yy}}^{1/2}\boldsymbol{d}_1$$

9

# Canonical Correlation Analysis (4)

- The (1st) canonical correlation is the largest singular value $\sigma_1$, where $\boldsymbol{S} = diag(\sigma_1, \sigma_2, \ldots)$.

- Similarly, the 2nd and higher canonical correlations are obtained from the other singular values.

- CCA has many interesting properties and close connections to PCA, PLS (partial least squares), CA (correspondence analysis), etc.
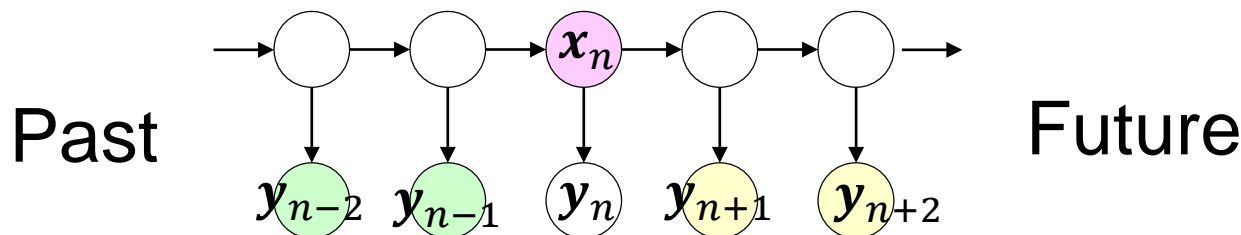
# Canonical Correlation Analysis of Time-series [Akaike 74][Akaike 75]

- Time-series of observation: $\{\cdots, \boldsymbol{y}_{n-1}, \boldsymbol{y}_n, \boldsymbol{y}_{n+1}, \cdots\}$

- "State" space representation of the system:

$$\boldsymbol{x}_{n+1} = F\boldsymbol{x}_n + \boldsymbol{w}_n$$
$$\boldsymbol{y}_n = H\boldsymbol{x}_n$$

- The "state" $\boldsymbol{x}_n$ can be interpreted in two ways:

  - Random variables that contain full information of the future to be expressed by the present and past

  - Random variables that contain full information of the past to be expressed by the present and future

Past  Future

11

# Canonical Correlation Analysis of Time-Series (cont.)

- "State" vector $x_n$ is the <span style="color:red">information interface between future</span> $y_n, y_{n+1}, \cdots$ <span style="color:red">and past</span> $y_{n-1}, y_{n-2}, \cdots$

- Predict future by past :

$$y_{n-1}, y_{n-2}, \cdots \quad \longrightarrow \quad u_n \quad \longrightarrow \quad y_n, y_{n+1}, \cdots$$

        Past                                      Future

- Postdict past by future :

$$y_n, y_{n+1}, \cdots \quad \longrightarrow \quad v_n \quad \longrightarrow \quad y_{n-1}, y_{n-2}, \cdots$$

        Future                                   Past

- Canonical vectors $u_n$ and $v_n$ should be maximally correlated

# Extension of CCA (1) : Kernel CCA

- Most of classical linear multivariate analysis methods can be non-linearized by (RKHS) kernel
  - Linear regression -> Kernel regression, SVR, GPR
  - Linear classifier -> SVM
  - Linear PCA -> Kernel PCA

Original space  Feature space

$x$   $\Phi(x)$

- Kernel CCA [Akaho 01][Bach 02]

Linear CCA

$$u = \boldsymbol{a}^T \boldsymbol{x}$$
$$v = \boldsymbol{b}^T \boldsymbol{y}$$

$$\max_{\boldsymbol{a},\boldsymbol{b}} \mathrm{cor}(u, v)$$

Non-linear mapping

$$\boldsymbol{x}^\Phi = \Phi(\boldsymbol{x})$$
$$\boldsymbol{y}^\Psi = \Psi(\boldsymbol{y})$$
$$\boldsymbol{a}^\Phi = \sum_i \alpha_i \cdot \Phi(\boldsymbol{x}_i)$$
$$\boldsymbol{b}^\Psi = \sum_i \beta_i \cdot \Psi(\boldsymbol{y}_i)$$

Kernel CCA

$$u = \boldsymbol{a}^{\Phi T} \boldsymbol{x}^\Phi$$
$$v = \boldsymbol{b}^{\Psi T} \boldsymbol{y}^\Psi$$

$$\max_{\boldsymbol{\alpha},\boldsymbol{\beta}} \mathrm{cor}(u, v)$$

13

# Learning Non-linear Dynamical Systems by Kernel CCA

- Y. Kawahara, T. Yairi, and K. Machida, "A kernel subspace method by stochastic realization for learning nonlinear dynamical systems", NIPS-2006
- Kernel CCA between past and future data
- A pioneering work of spectral learning of dynamical systems in machine learning community
- But, it was too early ..

Non-linear dynamical system

Linear dynamical system in kernel feature space

$$x(t+1) = g(x(t), u(t)) + v$$
$$y(t) = h(x(t), u(t)) + w,$$

$$x(t+1) = A^\phi x(t) + B^\phi \phi_u(u(t)) + K^\phi e(t),$$
$$\phi_y(y(t)) = C^\phi x(t) + D^\phi \phi_u(u(t)) + e(t),$$

14

# Extension of CCA (2): Probabilistic CCA

- Probabilistic canonical correlation analysis (PCCA) [Bach & Jordan 06]
  - Probabilistic (latent variable) model
  - Similar to Probabilistic PCA, variational autoencoder, etc.

PPCA, VAE          PCCA          Mixture of PCCA

Latent vector
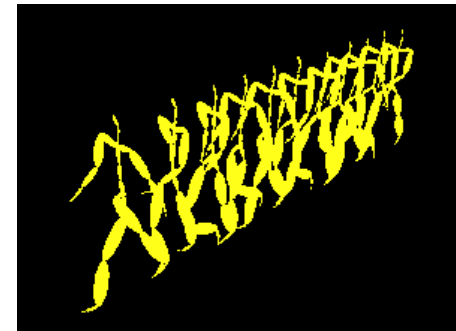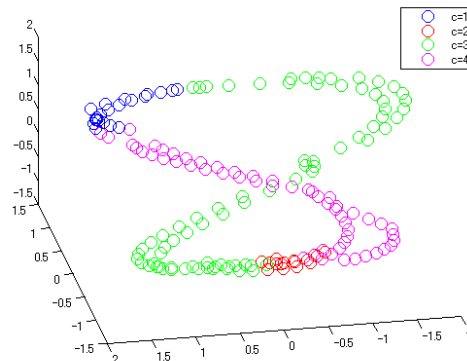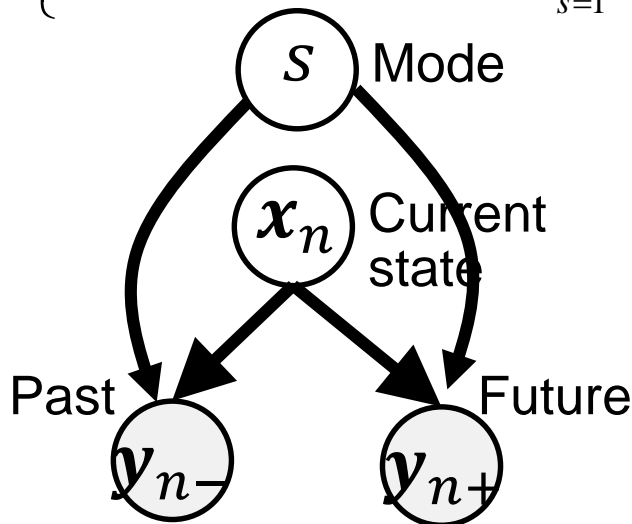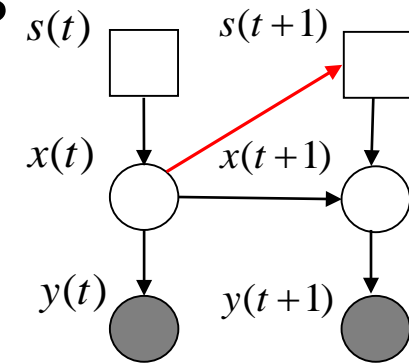(Low-dim.)

Observation
(High-dim.)

# Learning Mixture of Linear Dynamical Systems

Masao Joko, Yoshinobu Kawahara, Takehisa Yairi, "Learning Non-linear Dynamical Systems by Alignment of Local Linear Models", 20th International Conference on Pattern Recognition (ICPR), pp. 1084-1087, 2010

## Mixture of locally linear dynamical systems

$$p(x(t+1) \mid x(t)) = \sum_{s=1}^{C} p(s(t+1) \mid x(t)) \, p(x(t+1) \mid x(t), s(t+1))$$

$$p(y(t) \mid x(t)) = \sum_{s=1}^{C} p(s(t) \mid x(t)) \, p(y(t) \mid x(t), s(t))$$

$s(t)$  $s(t+1)$

$x(t)$  $x(t+1)$

$y(t)$  $y(t+1)$

$s$ Mode

$x_n$ Current state

Past $y_{n-}$    Future $y_{n+}$

c=1
c=2
c=3
c=4

# 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.)

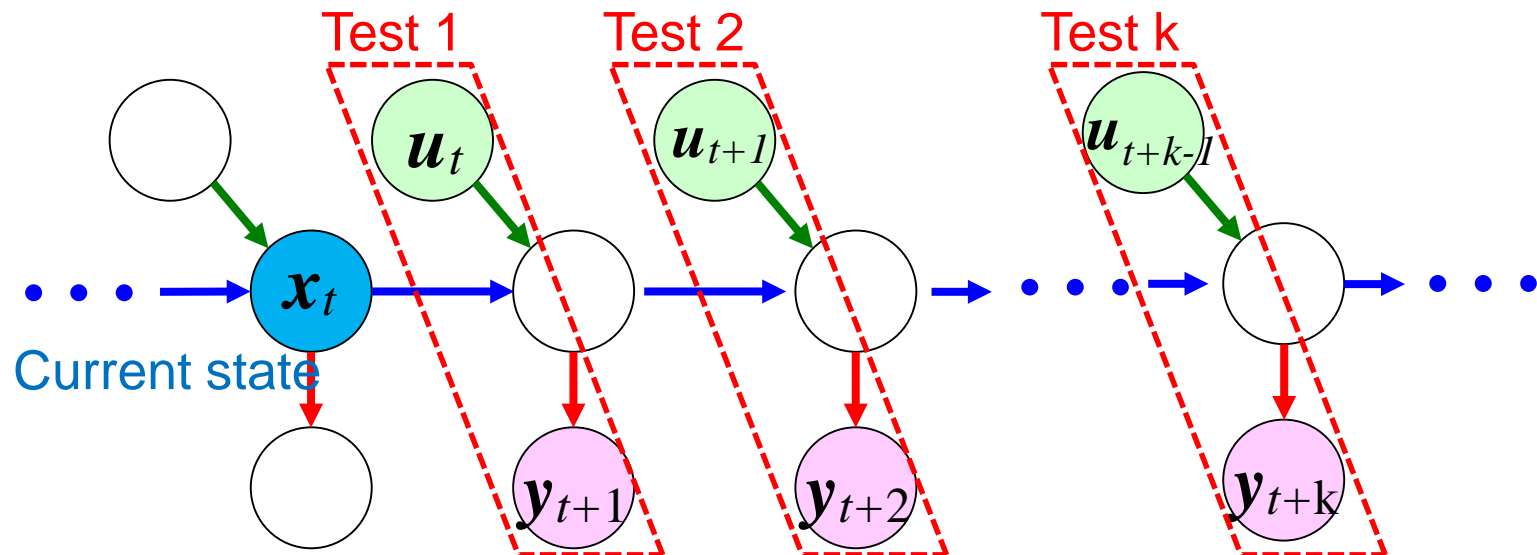The seminal work of [Hsu 09] was rapidly extended

- – 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] Extention 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]

# 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 ≈ Guessing the current state

# 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]

# Epilogue

- We focused on the methods of "learning dynamical systems" or <span style="color:red">machine learning approaches to system identification</span> problem
  - Learning a <span style="color:red">state space model</span> from observation data
- They are roughly classified into two approaches
  1. Maximum-likelihood approach
  2. Spectral learning approach
- Topics not covered this time:
  - (Deep) neural networks for time-series (RNN, etc.)
  - Koopeman operator, dynamic mode decomposition

# Final Task (of Yairi's Part)

- 適当な多変量時系列データに対して、システム同定・機械学習（動的システム学習）の手法を適用することによって、そのデータの発生源であるシステムの挙動モデルを推定し、その結果を考察せよ。用いた手法のアルゴリズム、データの出典や参考文献、等を明記すること。

- 締切： 2018年8月10日(金)

- 9月修了予定者は別途相談すること。

- 提出方法：ITC-LMS で提出

- 注：堀先生の課題も必ずやること