

Subspace Identification and Spectral Learning of Dynamical Systems

Jun.21, 2018

Takehisa YAIRI (矢入健久)

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

Outline

- What is subspace identification ?
- Least square linear regression and singular value decomposition
- Deterministic and stochastic subspace identification
 - Extended state space representation
 - Canonical correlation analysis
- Spectral learning of dynamical systems

(Review) Learning LDS

- Model:
$$\begin{cases} \mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \mathbf{v} & \mathbf{v} \sim N(\mathbf{0}, \mathbf{Q}) \\ \mathbf{y}_t = \mathbf{C}\mathbf{x}_t + \mathbf{w} & \mathbf{w} \sim N(\mathbf{0}, \mathbf{R}) \end{cases}$$
- Given :
 - Observation sequence : $\mathbf{y}_{1:T}$
 - Input sequence : $\mathbf{u}_{1:T}$
- Find :
 - System matrices : $\mathbf{A}, \mathbf{B}, \mathbf{C}$,
 - Noise covariance matrices: \mathbf{Q}, \mathbf{R}
 - Initial state distribution: $p(\mathbf{x}_1) \sim N(\mathbf{m}_0, \mathbf{V}_0)$
 - State sequence: $\mathbf{x}_{1:T}$

(Review) EM Algorithm for LDS

- Initialize model parameters
 - $\Theta^{(0)} = \{A^{(0)}, B^{(0)}, C^{(0)}, Q^{(0)}, R^{(0)}, m_0^{(0)}, V_0^{(0)}\}$
- Repeat until convergence:
 - [E-step] Inference by Kalman (RTS) Smoothing
 - Filtering : Compute $p(x_t | y_{1:t}, u_{1:t}, \Theta^{(t)}) = N(x_t | m_t, V_t)$
 - Smoothing: Compute $p(x_t | y_{1:T}, u_{1:T}, \Theta^{(t)}) = N(x_t | \hat{m}_t, \hat{V}_t)$
 - Cross-covariance of x_t and x_{t+1} is also necessary
 - [M-step] Model update
 - $\Theta^{(t+1)} \leftarrow \underset{\Theta}{\operatorname{argmax}} E_{x_{1:T}} [\log p(x_{1:T}, y_{1:T} | u_{1:T}, \Theta)]$

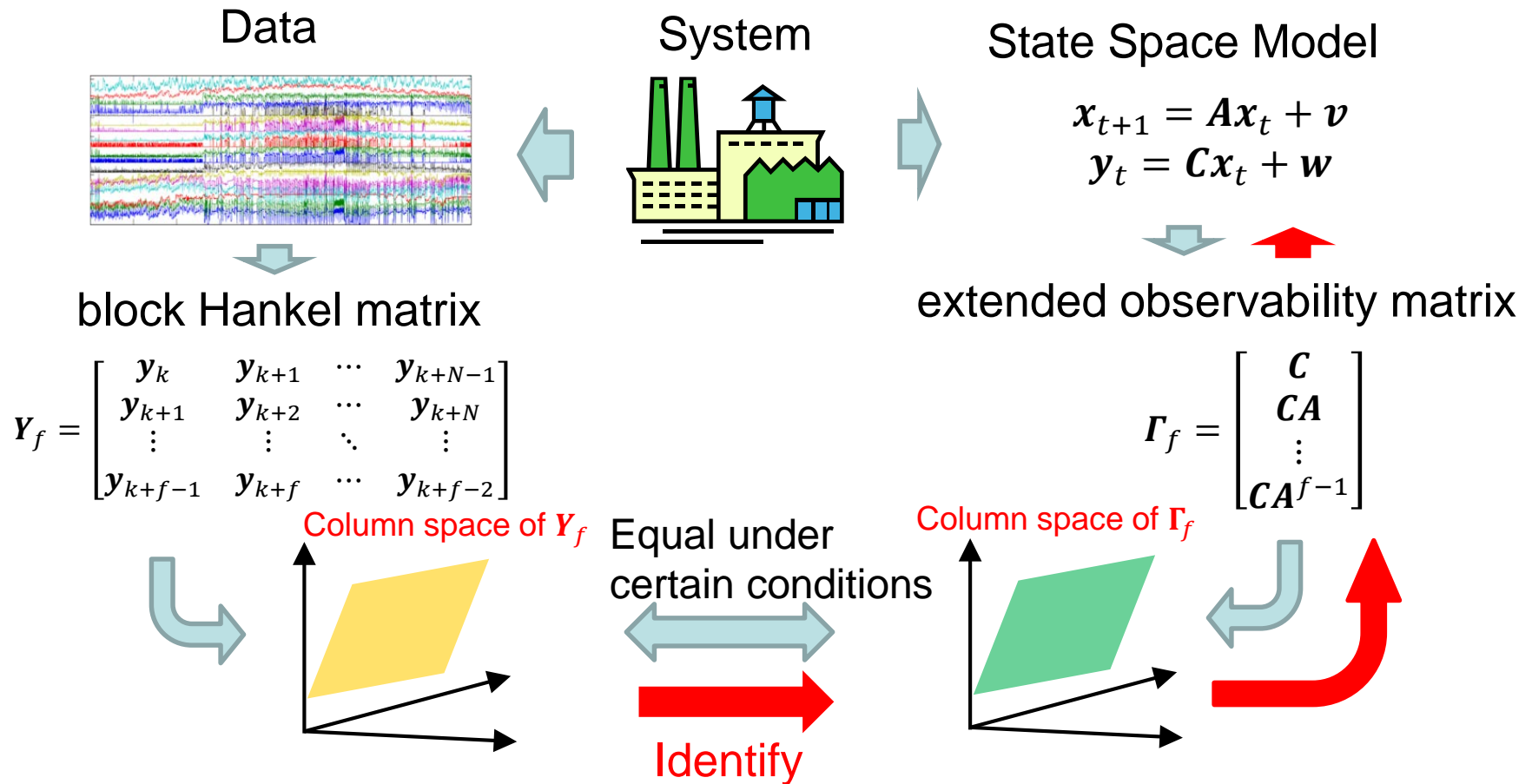
Pros:

- Exact (for LDS), highly applicable (to complicated models)

Cons:

- Prone to local optima, many iterations (slow)

What is "subspace" identification method ?



- **Advantages:**

- Global optimum is obtained without iteration by linear algebraic operations

Least Squares Linear Regression, and Singular Value Decomposition

Linear Regression

- Linear regression model (with vector output)

$$\mathbf{y} = \mathbf{\Gamma} \mathbf{x} + \mathbf{v}$$

Output (target) variable (vectors) Input (exploratory) variable error (residual)

- Collected data

$$\underbrace{[\mathbf{y}_1 \quad \mathbf{y}_2 \quad \cdots \quad \mathbf{y}_N]}_{\mathbf{Y}} = \mathbf{\Gamma} \underbrace{[\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_N]}_{\mathbf{X}} + \mathbf{V}$$



$$\mathbf{Y} = \mathbf{\Gamma} \mathbf{X} + \mathbf{V}$$



Estimate

Least Squares Solution

- Minimize sum of squares of residuals:

$$Loss = \|V\|_F^2 = \|Y - \Gamma X\|_F^2$$

$\|\cdot\|_F^2$: Frobenius norm

$$= Tr[(Y - \Gamma X)^T (Y - \Gamma X)]$$

$$\frac{\partial}{\partial \Gamma} Loss = -2YX^T + 2\Gamma XX^T = \mathbf{0}$$

Gradient is zero
at the minimum

$$\hat{\Gamma} = Y \underbrace{X^T (XX^T)^{-1}}_{\text{Pseudo inverse of } X}$$

Least square solution

$$\hat{Y} = \hat{\Gamma} X = Y X^T (XX^T)^{-1} X$$

Prediction

$$V = Y - \hat{Y} = Y \left(I - X^T (XX^T)^{-1} X \right)$$

Residual

Geometric Interpretation of Least Squares

Prediction

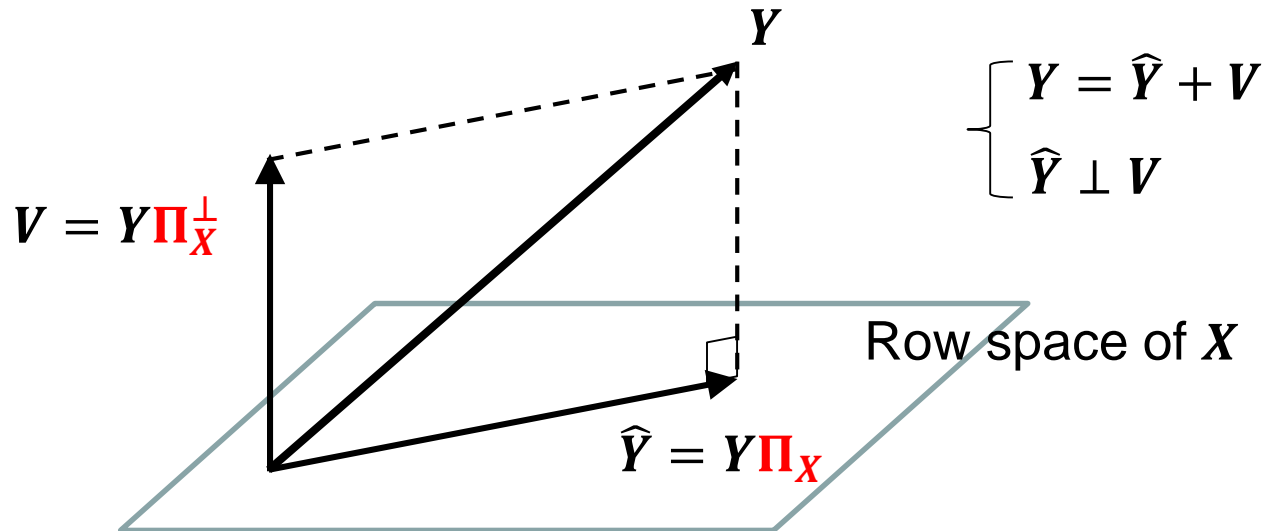
$$\hat{Y} = \hat{\Gamma}X = YX^T(XX^T)^{-1}X \\ \equiv Y\Pi_X$$

Orthogonal projection of Y onto the space spanned by the rows of X (row space)

Redidual

$$V = Y - \hat{Y} = Y(I - X^T(XX^T)^{-1}X) \\ = Y(I - \Pi_X) \equiv Y\Pi_X^\perp$$

Orthogonal projection of Y onto the complement of row space of X



When Two Sets of Explanatory Variables are given

- Consider there are two sets of inputs X and U :

$$Y = \Gamma X + HU + V$$

Unknown Unknown
 ↙ ↘
 ↗ ↘
 Known (data) Residual

- Assume V is independent of both X and U
- Multiply Π_U^\perp from the right side

$$Y\Pi_U^\perp = \Gamma X\Pi_U^\perp + \underbrace{HU\Pi_U^\perp}_0 + \underbrace{V\Pi_U^\perp}_V$$

$$\Rightarrow Y\Pi_U^\perp = \Gamma X\Pi_U^\perp + V$$

$$\Rightarrow \hat{\Gamma} = Y\Pi_U^\perp X^T (X\Pi_U^\perp X^T)^{-1}$$

Remarks:

- \hat{H} can also be similarly obtained
- Known as oblique projection

What if X is not given ?

Consider both Γ and X are unknown

$$Y = \Gamma X + W$$

But assume the number of columns of Γ is known

The diagram shows a 4x4 grid labeled Y followed by an approximation symbol \approx . To the right is a 4x3 grid labeled Γ with a bracket above it indicating d columns. This is followed by a multiplication symbol \times and a 3x2 grid labeled X .

This is the problem of **low-rank approximation** of a matrix



(Truncated) Singular Value Decomposition (SVD)

$$Y = U_d S_d V_d^T + W$$

$$\text{where, } S_d = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_d \end{bmatrix}$$

$$\left\{ \begin{array}{l} \hat{\Gamma} = U_d S_d^{1/2} \\ \hat{X} = S_d^{1/2} V_d^T \end{array} \right.$$

(Note that solution is not unique)

Innovation Form of State Space Model

- Get back to the state space model

$$\begin{cases} \mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{v}_k \\ \mathbf{y}_k = \mathbf{C}\mathbf{x}_k + \mathbf{w}_k \end{cases}$$

- Kalman filter representation

$$\hat{\mathbf{x}}_{k+1} = \mathbf{A}\hat{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k + \mathbf{K}(\mathbf{y}_k - \mathbf{C}\hat{\mathbf{x}}_k)$$

- Define the *innovation* : $\mathbf{e}_k = \mathbf{y}_k - \mathbf{C}\hat{\mathbf{x}}_k$

$$\begin{cases} \hat{\mathbf{x}}_{k+1} = \mathbf{A}\hat{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k + \mathbf{K}\mathbf{e}_k \\ \mathbf{y}_k = \mathbf{C}\hat{\mathbf{x}}_k + \mathbf{e}_k \end{cases}$$

Note that inputs \mathbf{u}_k and outputs \mathbf{y}_k are unchanged

Future Outputs

Innovation form

$$\begin{cases} \hat{\mathbf{x}}_{k+1} = \mathbf{A}\hat{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k + \mathbf{K}\mathbf{e}_k \\ \mathbf{y}_k = \mathbf{C}\hat{\mathbf{x}}_k + \mathbf{e}_k \end{cases}$$

1-step
ahead

$$\mathbf{y}_{k+1} = \mathbf{C}\hat{\mathbf{x}}_{k+1} + \mathbf{e}_{k+1} = \mathbf{C}\mathbf{A}\hat{\mathbf{x}}_k + \mathbf{C}\mathbf{B}\mathbf{u}_k + \mathbf{C}\mathbf{K}\mathbf{e}_k + \mathbf{e}_{k+1}$$

⋮

⋮

j -step
ahead

$$\mathbf{y}_{k+j} = \mathbf{C}\mathbf{A}^j\hat{\mathbf{x}}_k + [\mathbf{C}\mathbf{A}^{j-1}\mathbf{B} \quad \dots \quad \mathbf{C}\mathbf{B} \quad \mathbf{0}] \begin{bmatrix} \mathbf{u}_k \\ \vdots \\ \mathbf{u}_{k+j} \end{bmatrix}$$

$$+ [\mathbf{C}\mathbf{A}^{j-1}\mathbf{K} \quad \dots \quad \mathbf{C}\mathbf{K} \quad \mathbf{I}] \begin{bmatrix} \mathbf{e}_k \\ \vdots \\ \mathbf{e}_{k+j} \end{bmatrix}$$



stack vertically

Extended
observability
matrix

$$\begin{aligned}
 \begin{bmatrix} \mathbf{y}_k \\ \mathbf{y}_{k+1} \\ \vdots \\ \mathbf{y}_{k+f-1} \end{bmatrix} &= \overbrace{\begin{bmatrix} \mathbf{C} \\ \mathbf{CA} \\ \vdots \\ \mathbf{CA}^{f-1} \end{bmatrix}}^{\mathbf{\Gamma}_f} \hat{\mathbf{x}}_k + \overbrace{\begin{bmatrix} \mathbf{0} & \mathbf{0} & \vdots & \mathbf{CA}^{f-2}\mathbf{B} \\ \mathbf{CB} & \mathbf{0} & \mathbf{CB} & \dots \\ \vdots & \mathbf{CB} & \ddots & \mathbf{CB} \\ \mathbf{CA}^{f-2}\mathbf{B} & \dots & \mathbf{CB} & \mathbf{0} \end{bmatrix}}^{\mathbf{H}_f} \begin{bmatrix} \mathbf{u}_k \\ \mathbf{u}_{k+1} \\ \vdots \\ \mathbf{u}_{k+f-1} \end{bmatrix} \\
 &+ \underbrace{\begin{bmatrix} \mathbf{I} & \mathbf{I} & \vdots & \mathbf{CA}^{f-2}\mathbf{K} \\ \mathbf{CK} & \mathbf{I} & \mathbf{CK} & \dots \\ \vdots & \mathbf{CK} & \ddots & \mathbf{CK} \\ \mathbf{CA}^{f-2}\mathbf{K} & \dots & \mathbf{CK} & \mathbf{I} \end{bmatrix}}_{\mathbf{G}_f} \begin{bmatrix} \mathbf{e}_k \\ \mathbf{e}_{k+1} \\ \vdots \\ \mathbf{e}_{k+f-1} \end{bmatrix}
 \end{aligned}$$

Relation between the current
state (estimate) $\hat{\mathbf{x}}_k$ and future
measurements $\mathbf{y}_{k:k+f-1}$, given
future inputs $\mathbf{u}_{k:k+f-1}$

\mathbf{H}_f and \mathbf{G}_f are Toeplitz matrices

 stack horizontally

Future output data
matrix

\mathbf{Y}_f

$$\begin{bmatrix} \mathbf{y}_k & \cdots & \mathbf{y}_{k+N-1} \\ \vdots & \cdots & \vdots \\ \mathbf{y}_{k+f-1} & \cdots & \mathbf{y}_{k+f-2} \end{bmatrix}$$

Kalman state
sequence

\mathbf{X}_k

$$= \mathbf{\Gamma}_f [\hat{\mathbf{x}}_k \quad \cdots \quad \hat{\mathbf{x}}_{k+N-1}] + \mathbf{H}_f$$

Future input data
matrix

\mathbf{U}_f

$$\begin{bmatrix} \mathbf{u}_k & \cdots & \mathbf{u}_{k+N-1} \\ \vdots & \cdots & \vdots \\ \mathbf{u}_{k+f-1} & \cdots & \mathbf{u}_{k+f-2} \end{bmatrix}$$

Given

Given

Weighted
residuals

$$\mathbf{Y}_f = \underbrace{\mathbf{\Gamma}_f \mathbf{X}_k}_{\text{Unknown}} + \mathbf{H}_f \mathbf{U}_f + \mathbf{G}_f \mathbf{E}_f$$

Unknown

$$+ \mathbf{G}_f \begin{bmatrix} \mathbf{e}_k & \cdots & \mathbf{e}_{k+N-1} \\ \vdots & \cdots & \vdots \\ \mathbf{e}_{k+f-1} & \cdots & \mathbf{e}_{k+f-2} \end{bmatrix}$$

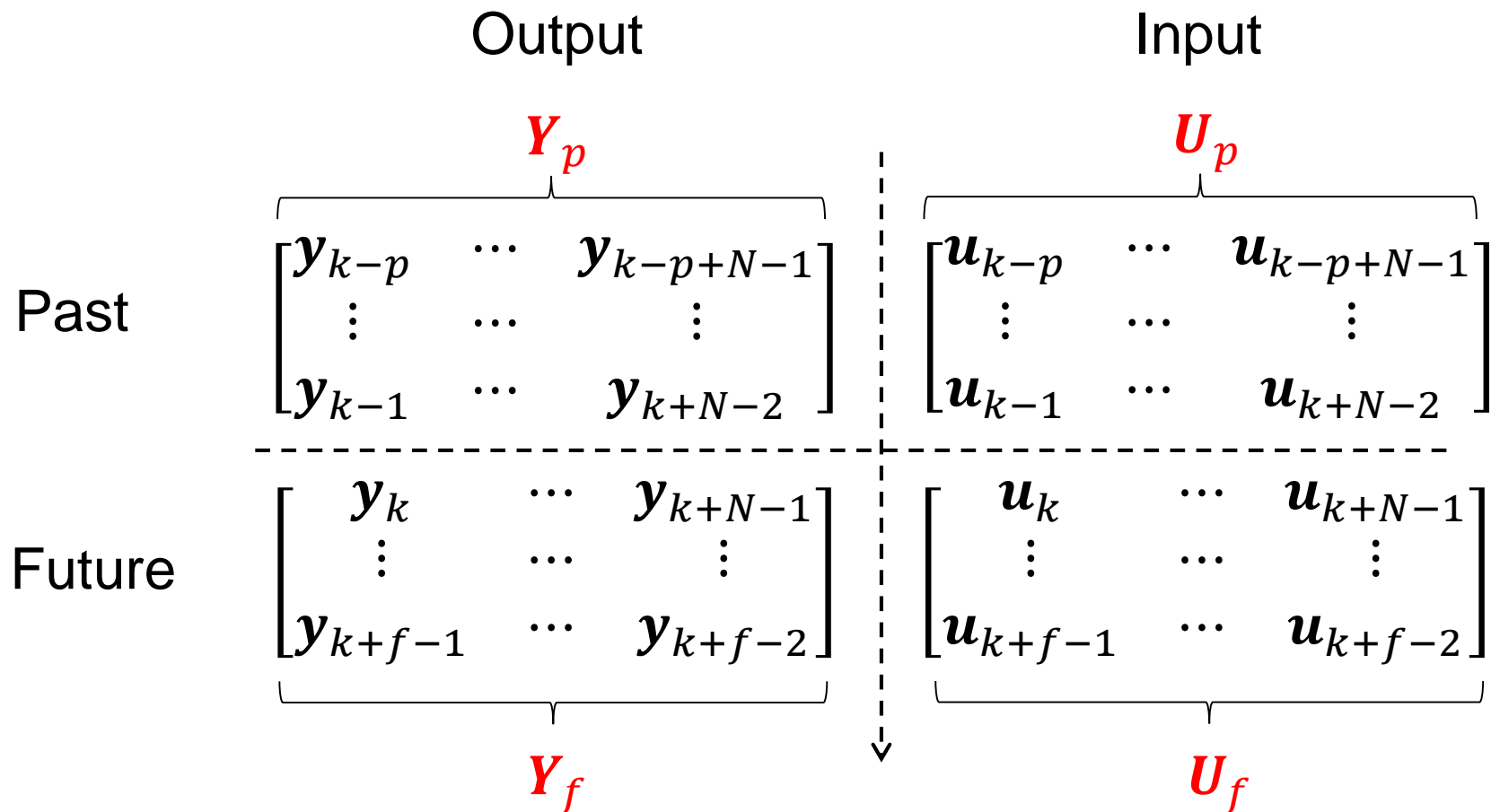
\mathbf{E}_f

Future noise data
matrix

$\mathbf{Y}_f, \mathbf{U}_f, \mathbf{E}_f$ are block Hankel matrices

Past and Future Data Matrices

Past input / output data matrices are also defined



Kalman State Sequence

- Kalman state (sequence) \mathbf{X}_k is unknown
- But, it is known it **can be estimated from past inputs and outputs !**

$$\mathbf{X}_k = \mathbf{L}_u \mathbf{U}_p + \mathbf{L}_y \mathbf{Y}_p = [\mathbf{L}_u \quad \mathbf{L}_y] \begin{bmatrix} \mathbf{U}_p \\ \mathbf{Y}_p \end{bmatrix} \equiv \mathbf{L}_z \mathbf{Z}_p$$

↑
Past inputs
and outputs

$$\mathbf{Y}_f = \mathbf{\Gamma}_f \mathbf{X}_k + \mathbf{H}_f \mathbf{U}_f + \mathbf{G}_f \mathbf{E}_f$$



$$\mathbf{Y}_f = \mathbf{\Gamma}_f \mathbf{L}_z \mathbf{Z}_p + \mathbf{H}_f \mathbf{U}_f + \mathbf{G}_f \mathbf{E}_f$$

Open Loop Assumption

- We assume the system is open loop
- Future innovation is uncorrelated to future inputs

$$\mathbf{E}_f \mathbf{U}_f^T = \mathbf{0} \Leftrightarrow \mathbf{E}_f \mathbf{\Pi}_{\mathbf{U}_f}^\perp = \mathbf{E}_f$$

- Future innovation is uncorrelated to past inputs and outputs

$$\mathbf{E}_f \mathbf{Z}_p^T = \mathbf{0}$$

Subspace Identification 1

(N4SID approach) [Qin 2004][Qin 2006]

- Data matrices Y_f , U_f and Z_p are given

$$Y_f = \Gamma_f L_z Z_p + H_f U_f + G_f E_f$$

- Multiply $\Pi_{U_f}^\perp$ from the right side

$$\begin{aligned} Y_f \Pi_{U_f}^\perp &= \Gamma_f L_z Z_p \Pi_{U_f}^\perp + H_f \overbrace{U_f \Pi_{U_f}^\perp}^{\text{Zero}} + G_f \overbrace{E_f \Pi_{U_f}^\perp}^{E_f} \\ &= \Gamma_f L_z Z_p \Pi_{U_f}^\perp + G_f E_f \end{aligned}$$

- Multiply Z_p^T from the right side

$$\begin{aligned} \underbrace{Y_f \Pi_{U_f}^\perp Z_p^T}_{\substack{\text{Known} \\ \text{(from data)}}} &= \Gamma_f L_z Z_p \Pi_{U_f}^\perp Z_p^T + G_f \overbrace{E_f Z_p^T}^{\text{Zero}} \\ &= \Gamma_f L_z Z_p \Pi_{U_f}^\perp Z_p^T \end{aligned}$$

N4SID Approach (Continued)

- Perform (truncated) SVD on $Y_f \Pi_{U_f}^\perp Z_p^T$

$$Y_f \Pi_{U_f}^\perp Z_p^T \approx \mathbf{U}_d \mathbf{S}_d \mathbf{V}_d^T \left(\approx \mathbf{\Gamma}_f L_z Z_p \Pi_{U_f}^\perp Z_p^T \right)$$

- Determine the extended observability matrix as:

$$\hat{\mathbf{\Gamma}}_f \leftarrow \mathbf{U}_d \mathbf{S}_d^{1/2}$$

$$\mathbf{\Gamma}_f = \begin{bmatrix} \mathbf{C} \\ \mathbf{CA} \\ \vdots \\ \mathbf{CA}^{f-1} \end{bmatrix}$$

⇒ System matrices A , C are obtained

Subspace Identification 2 : Regression Approach

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

$$\underbrace{Y_f \Pi_{U_f}^\perp}_{\text{Known}} = \underbrace{\Gamma_f L_z}_{\text{Unknown}} \underbrace{Z_p \Pi_{U_f}^\perp}_{\text{Known}} + \underbrace{G_f E_f}_{\text{Residual}}$$

Target variables Explanatory variables

Same as the first approach
 Smaller is better

Obtained by (multi-variate) linear regression:

$$\begin{aligned} \widehat{\Gamma_f L_z} &= Y_f \Pi_{U_f}^\perp \left(\Pi_{U_f}^\perp Z_p^T \right) \left(Z_p \Pi_{U_f}^\perp \Pi_{U_f}^\perp Z_p^T \right)^{-1} \\ &= Y_f \Pi_{U_f}^\perp Z_p^T \left(Z_p \Pi_{U_f}^\perp Z_p^T \right)^{-1} \approx \mathbf{U}_d \mathbf{S}_d \mathbf{V}_d^T \quad (\text{SVD}) \end{aligned}$$

$$\hat{\Gamma}_f \leftarrow \mathbf{U}_d \mathbf{S}_d^{1/2}$$

Subspace Identification 3 : CCA Approach*

$$\underbrace{Y_f \Pi_{U_f}^\perp}_{\text{Known}} = \underbrace{\Gamma_f}_{\text{Unknown}} \underbrace{L_z Z_p \Pi_{U_f}^\perp}_{\text{Known}} + \underbrace{G_f E_f}_{\text{Residual}}$$

Variable set 1 Variable set 2

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

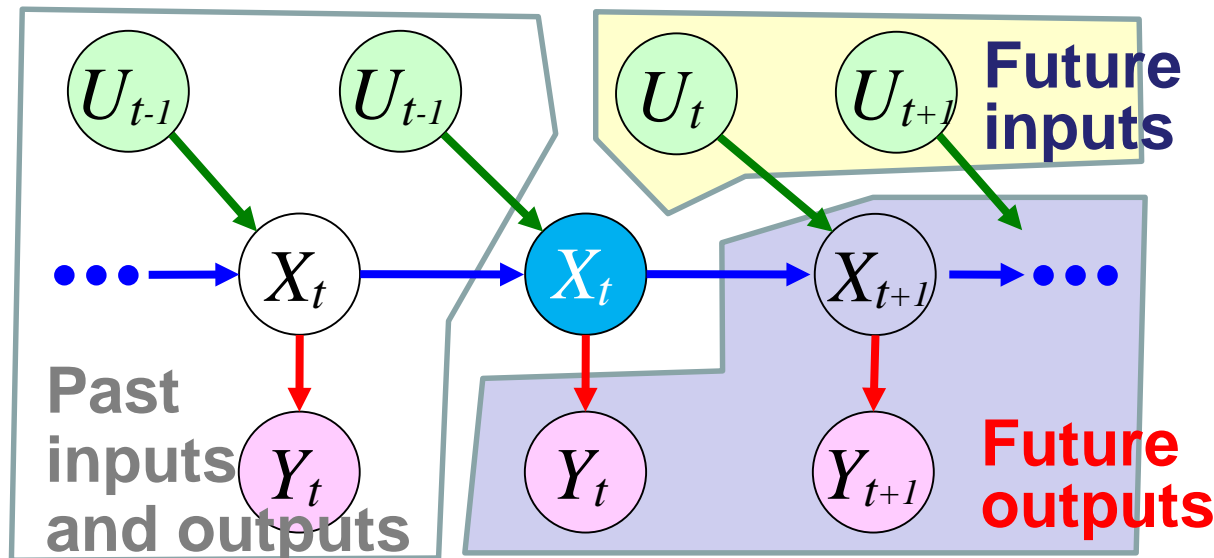
(*) Also known as canonical variate analysis (CVA)

CCA-based Subspace Identification (Continued)

- Additionally, state vectors are obtained by

$$\hat{\mathbf{X}}_k = \mathbf{S}_d^{1/2} \mathbf{V}_d^T \mathbf{W}_c \mathbf{Z}_p$$

- This result is very interesting, because the "state" is obtained by **CCA between "past" and "future"** data sets.
- This idea goes back to [Akaike 1975].



"The state vector carries information necessary to predict the future output based on the past."

[Katayama 2005]

Reference

- [Qin 2006] S. Joe Qin, "An overview of subspace identification", Computers & Chemical Engineering, Volume 30, Issues 10–12, 2006, Pages 1502-1513
 - Presentation file [Qin 2004] is also available (<http://people.duke.edu/~hpgavin/SystemID/References/Qin-SubspaceID-2004.pdf>)
- [Overschee 1996] Peter VAN OVERSCHEE & Bart DE MOOR, "SUBSPACE IDENTIFICATION FOR LINEAR SYSTEMS - Theory - Implementation - Applications", Springer, 1996
- [Katayama 2005] Toru Katayama, "Subspace Methods for System Identification ", Springer, 2005
- [Akaike 1975] H.Akaike, "Markovian Representation of Stochastic Processes by Canonical Variables", SIAM J. CONTROL, Vol. 13, No. 1, pp.162-173, January 1975

Appendix:

Canonical Correlation Analysis (CCA)

Canonical Correlation Analysis (1)

There are two random variables x and y

$$\mathbf{x} \in R^m \quad (\text{m-dim.}) \quad \mathbf{y} \in R^n \quad (\text{n-dim.})$$

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

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

Define covariance matrices of / between x and y

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

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

Canonical Correlation Analysis (2)

Consider to construct synthetic variable u and v by linear combinations of \mathbf{x} and \mathbf{y} , respectively.

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

$$v = \mathbf{b}^T \mathbf{y} = b_1 y_1 + b_2 y_2 + \cdots + b_n y_n$$

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

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

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

Canonical Correlation Analysis (3)

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

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

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

The goal of CCA is to solve:

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

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

Canonical Correlation Analysis (4)

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

$$\mathbf{c} = \Sigma_{xx}^{1/2} \mathbf{a} \quad , \text{ and } \mathbf{d} = \Sigma_{yy}^{1/2} \mathbf{b}$$

The problem turns to be

$$(\mathbf{c}_1, \mathbf{d}_1) = \underset{\|\mathbf{c}\|=1, \|\mathbf{d}\|=1}{\operatorname{argmax}} \quad \mathbf{c}^T \left(\Sigma_{xx}^{-T/2} \Sigma_{xy} \Sigma_{xx}^{-1/2} \right) \mathbf{d}$$

This problem can be solved by SVD !

$$\Sigma_{xx}^{-T/2} \Sigma_{xy} \Sigma_{xx}^{-1/2} = \mathbf{U} \mathbf{S} \mathbf{V}^T \quad \Rightarrow \quad \mathbf{c}_1 = \mathbf{u}_1, \mathbf{d}_1 = \mathbf{v}_1$$

$$\Rightarrow \quad \mathbf{a}_1 = \Sigma_{xx}^{1/2} \mathbf{u}_1, \text{ and } \mathbf{b}_1 = \Sigma_{yy}^{1/2} \mathbf{v}_1$$