

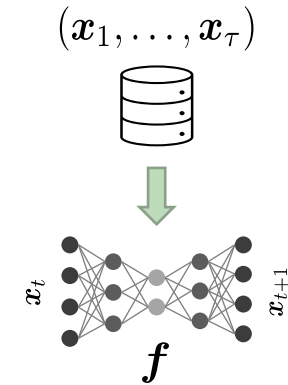
# Learning Multiple Nonlinear Dynamical Systems with Side Information

Naoya Takeishi (HES-SO/RIKEN), Yoshinobu Kawahara (KyushuU/RIKEN)

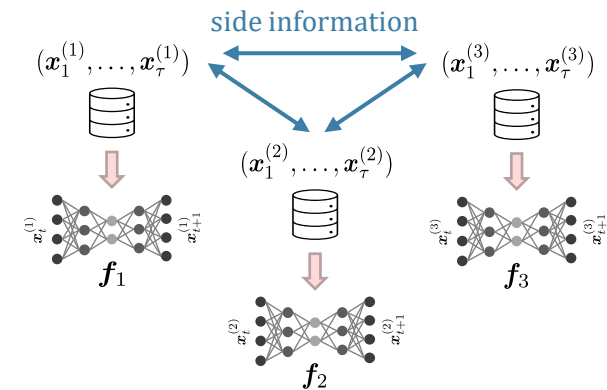
The 59th IEEE CDC, December 2020

# Problem Setting

- Machine learning of dynamical system
  - Given a dataset (i.e., a sequence of states)  $(\mathbf{x}_1, \dots, \mathbf{x}_\tau), \mathbf{x} \in \mathbb{R}^d$ ,
  - learn  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  s.t.  $\mathbf{x}_{t+1} \approx \mathbf{f}(\mathbf{x}_t)$
  - e.g., linear models, kernel machines, neural nets, ...

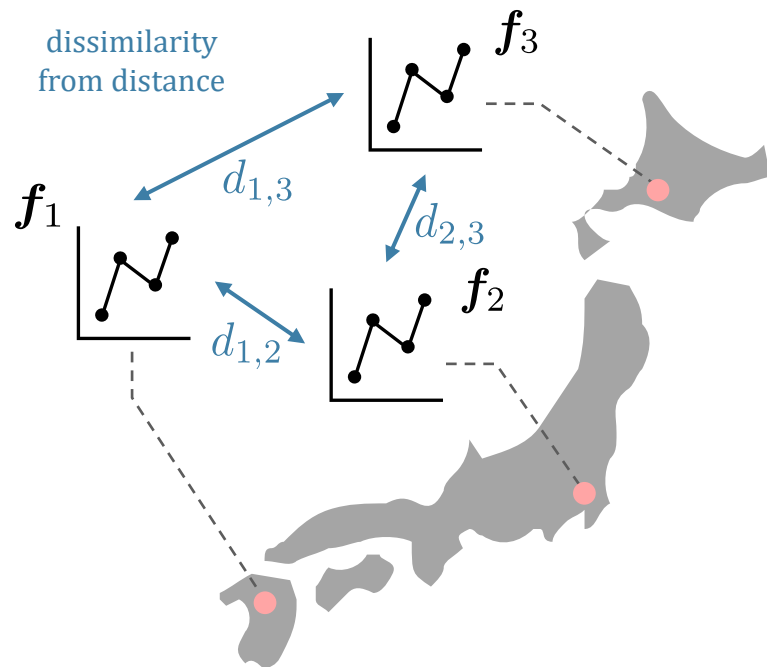


- Machine learning of **multiple** dynamical systems
  - Given multiple datasets  $(\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_{\tau_1}^{(1)}), \dots, (\mathbf{x}_1^{(n)}, \dots, \mathbf{x}_{\tau_n}^{(n)})$ ,
  - learn  $\mathbf{f}_1, \dots, \mathbf{f}_n$  s.t.  $\mathbf{x}_{t+1}^{(1)} \approx \mathbf{f}_1(\mathbf{x}_t^{(1)}), \dots, \mathbf{x}_{t+1}^{(n)} \approx \mathbf{f}_n(\mathbf{x}_t^{(n)})$
  - Side information on **relation between datasets/dynamics** is often available and helpful for learning

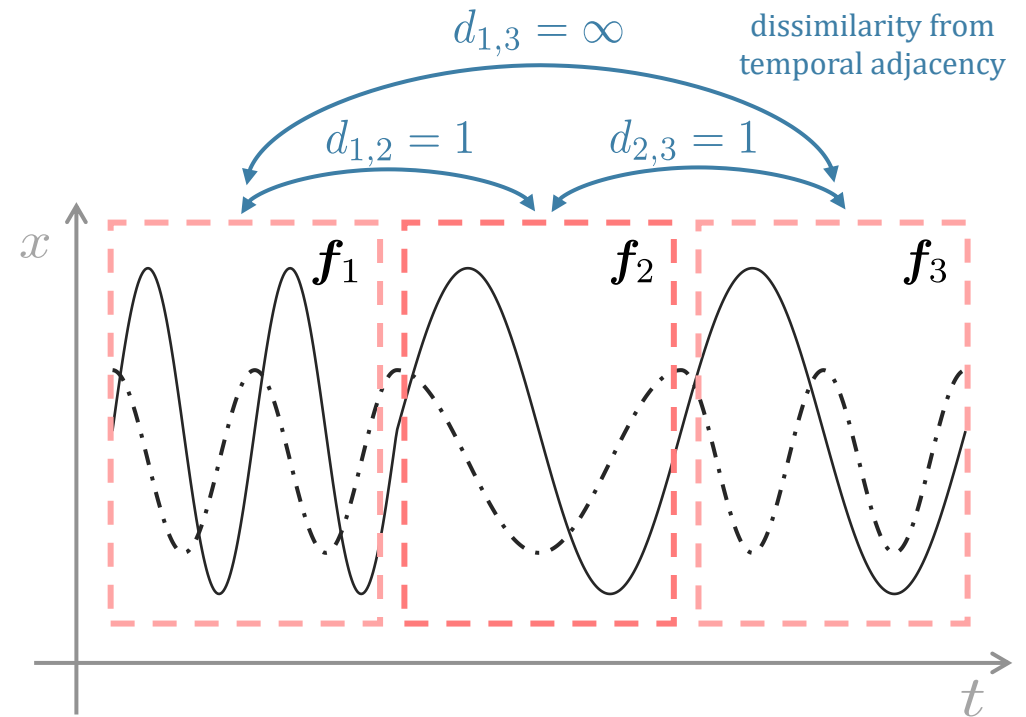


# Examples of Side Information

When learning dynamics from **geometrically distributed sensors**, **distance between sensors** may be informative for similarity of dynamics, e.g., **closer sensors measure similar dynamics**.

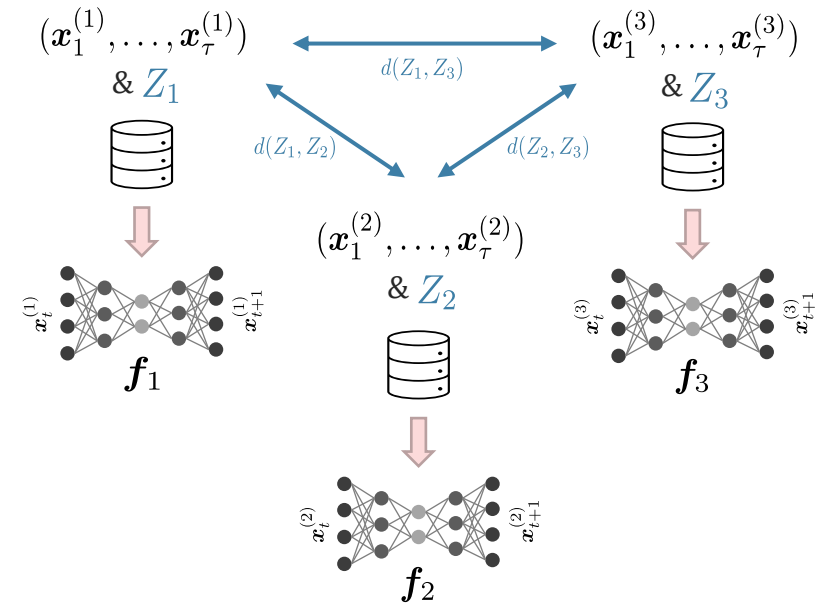


Learning **time-varying dynamical systems** can also be handled as a special case, where we may use  $t$  as **side information**, e.g., **dynamics in adjacent time periods are similar**.



# Problem Setting Again

- Input:
  - $n$  sets of measurements (state sequences)  
 $(\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_{\tau_1}^{(1)}), \dots, (\mathbf{x}_1^{(n)}, \dots, \mathbf{x}_{\tau_n}^{(n)})$
  - $n$  sets of side information  $Z_1, \dots, Z_n$ 
    - ✓ along with some dissimilarity measure  $d_Z(Z_i, Z_j)$
    - ✓  $Z$  may be location of sensors, timestamps, ..., additional measurements, text description, ...
- Output: dynamics models  $f_1, \dots, f_n$   
such that  $\mathbf{x}_{t+1}^{(1)} \approx f_1(\mathbf{x}_t^{(1)}), \dots, \mathbf{x}_{t+1}^{(n)} \approx f_n(\mathbf{x}_t^{(n)})$
- Core idea: Regularization with side information (next slide)
  - if  $Z_i$  and  $Z_j$  are similar (i.e.,  $d_Z(Z_i, Z_j)$  is small), then  $f_i$  and  $f_j$  should also be similar



# Regularization with Side Information

- Regularization with side information
  - if  $Z_i$  and  $Z_j$  are similar (i.e.,  $d_{\mathcal{Z}}(Z_i, Z_j)$  is small), then  $f_i$  and  $f_j$  should also be similar (i.e., some dissimilarity  $d_{\text{DS}}(f_i, f_j)$  is small)

- Formulation as multi-task learning for dynamical systems

$$\underset{\mathbf{f}_1, \dots, \mathbf{f}_n}{\text{minimize}} \underbrace{\frac{1}{n} \sum_{i=1}^n L_i(\mathbf{f}_i; \{\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{\tau_i}^{(i)}\})}_{\text{original objective function to learn } \mathbf{f}_1, \dots, \mathbf{f}_n} + \lambda \cdot \underbrace{\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n \frac{d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)}{d_{\mathcal{Z}}(Z_i, Z_j)}}_{\text{regularization term to make } d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) \propto d_{\mathcal{Z}}(Z_i, Z_j)}$$

e.g., squared loss  $L_i = \sum_{t=1}^{\tau_i-1} \|\mathbf{f}_i(\mathbf{x}_t^{(i)}) - \mathbf{x}_{t+1}^{(i)}\|_2^2$

- What  $d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)$  should be? (next slide)

# Dissimilarity Measure of Dynamical Systems

$$\underset{\mathbf{f}_1, \dots, \mathbf{f}_n}{\text{minimize}} \underbrace{\frac{1}{n} \sum_{i=1}^n L_i(\mathbf{f}_i; \{\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{\tau_i}^{(i)}\})}_{\text{original objective function to learn } \mathbf{f}_1, \dots, \mathbf{f}_n} + \lambda \cdot \underbrace{\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n \frac{d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)}{d_{\mathcal{Z}}(Z_i, Z_j)}}_{\text{regularization term to make } d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) \propto d_{\mathcal{Z}}(Z_i, Z_j)}$$

e.g., squared loss  $L_i = \sum_{t=1}^{\tau_i-1} \|\mathbf{f}_i(\mathbf{x}_t^{(i)}) - \mathbf{x}_{t+1}^{(i)}\|_2^2$

- $d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)$  should measure the dissimilarity between  $\mathbf{f}_i$  and  $\mathbf{f}_j$
- We use an operator-theoretic metric [Ishikawa+ 18]

$$d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) = \sqrt{1 - \frac{k_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)^2}{k_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_i)k_{\text{DS}}(\mathbf{f}_j, \mathbf{f}_j)}}, \quad k_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)^2 = \text{tr} \left( \bigwedge^m \sum_{t=0}^{\tau} (K_j^t)^* K_i^t \right)$$

$K_i, K_j$ : Perron-Frobenius operators in RKHS

- because it is applicable to nonlinear  $\mathbf{f}$  & agnostic of parametric form of  $\mathbf{f}$

# Dissimilarity Measure of Dynamical Systems

cont'd

$$d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) = \sqrt{1 - \frac{k_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)^2}{k_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_i)k_{\text{DS}}(\mathbf{f}_j, \mathbf{f}_j)}}, \quad k_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)^2 = \text{tr} \left( \bigwedge^m \sum_{t=0}^{\tau} (K_j^t)^* K_i^t \right)$$

- Let  $\mathcal{H}$  be an RKHS equipped with a kernel function  $k_{\mathcal{H}}(\cdot, \cdot)$ .
- Let  $\phi(\mathbf{x}) = k_{\mathcal{H}}(\mathbf{x}, \cdot)$  be the corresponding feature map.
- For a nonlinear dynamical system  $\mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t)$ , consider linear operator  $K$ :

$$K\phi(\mathbf{x}) = \phi(\mathbf{f}(\mathbf{x})),$$

which is called **Perron–Frobenius operator in RKHS** corresponding to  $\mathbf{f}$ .

- $K$  can be **estimated from trajectory**  $(\mathbf{x}_1, \dots, \mathbf{x}_\tau)$  with the kernel  $k_{\mathcal{H}}(\cdot, \cdot)$ .

$d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)$  is (almost everywhere) differentiable wrt. the parameters of  $\mathbf{f}_i$  &  $\mathbf{f}_j$

# Experiment: Datasets

- Synthetic datasets

- Van der Pol oscillator (→)

- ✓  $i$ -th dataset is generated via  $\ddot{x} - \mu_i(1 - x^2)\dot{x} + x = 0$

- ✓ Side information:  $d_Z(Z_i, Z_j) = |\mu_i - \mu_j|^2$

- Rössler system

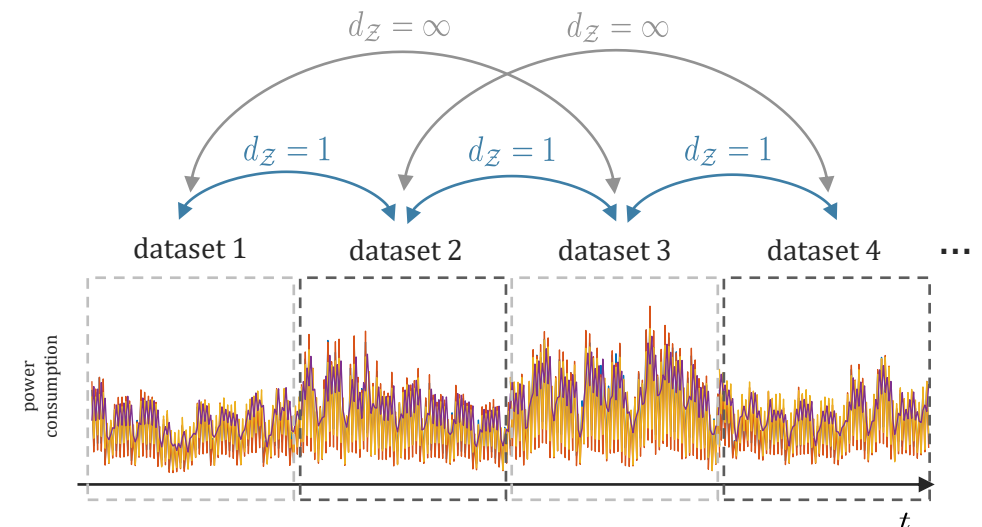
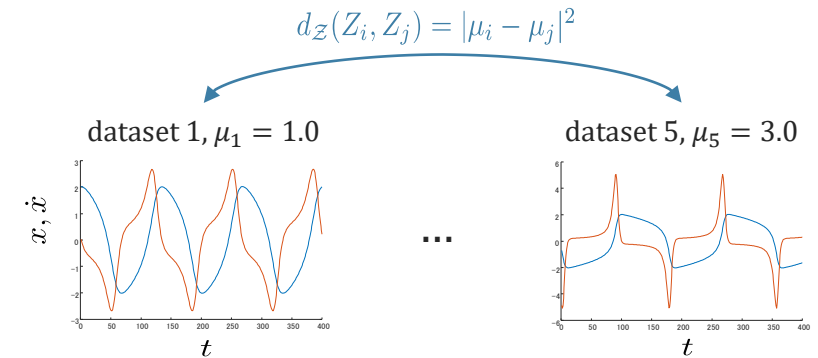
- Real-world datasets

- Solar power production time-series

- ✓ Side information:  $Z_i = \text{location of } i\text{-th plant}$

- Power consumption time-series (→)

- ✓ Side information:  $d_Z(Z_i, Z_j) = 1$  if  $j = i + 1$





# Experiment: Baseline Methods

1. No multi-task regularization

2. Multi-task w/ Parameter L2 dist.  $d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) = \sum_k |\theta_{i,k} - \theta_{j,k}|^2$

3. Multi-task w/ Parameter L1 dist.  $d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) = \sum_k |\theta_{i,k} - \theta_{j,k}|$

$\theta_{i,k}$  :  $k$ -th parameter of  $\mathbf{f}_i$

cf. Proposed method (= multi-task w/ dynamics dist.)

$$\underset{\mathbf{f}_1, \dots, \mathbf{f}_n}{\text{minimize}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n L_i(\mathbf{f}_i; \{\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{\tau_i}^{(i)}\})}_{\text{original objective function to learn } \mathbf{f}_1, \dots, \mathbf{f}_n} + \lambda \cdot \underbrace{\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n \frac{d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)}{d_{\mathcal{Z}}(Z_i, Z_j)}}_{\text{regularization term to make } d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) \propto d_{\mathcal{Z}}(Z_i, Z_j)}$$

original objective function to learn  $\mathbf{f}_1, \dots, \mathbf{f}_n$   
e.g., squared loss  $L_i = \sum_{t=1}^{\tau_i-1} \|\mathbf{f}_i(\mathbf{x}_t^{(i)}) - \mathbf{x}_{t+1}^{(i)}\|_2^2$

regularization term to make  $d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) \propto d_{\mathcal{Z}}(Z_i, Z_j)$

# Experiment: Results

Test one-step prediction errors (smaller is better)

	Baseline 1 No regularization	Baseline 2 Parameter L2 dist.	Baseline 3 Parameter L1 dist.	Proposed method (dynamics dist.)
VDP	$2.936 (.56) \times 10^{-1}$	$2.402 (.82) \times 10^{-1}$	$2.659 (.51) \times 10^{-1}$	<b><math>2.272 (.71) \times 10^{-1}</math></b>
RÖSSLER	$1.358 (.05) \times 10^{-2}$	$1.359 (.05) \times 10^{-2}$	$1.358 (.05) \times 10^{-2}$	<b><u><math>1.319 (.05) \times 10^{-2}</math></u></b>
SOLAR	$9.231 (.01) \times 10^{-4}$	<b><math>9.229 (.01) \times 10^{-4}</math></b>	<b><math>9.226 (.01) \times 10^{-4}</math></b>	<b><u><math>9.101 (.03) \times 10^{-4}</math></u></b>
DEMAND	$1.486 (.03) \times 10^{-3}$	$1.487 (.03) \times 10^{-3}$	$1.488 (.03) \times 10^{-3}$	<b><u><math>1.439 (.03) \times 10^{-3}</math></u></b>

Significant difference, **bold**: from (A) / underline: from (B), by paired  $t$ -test at  $p < 10^{-3}$ .

- Proposed regularization achieves better prediction
  - Probably, L2/L1 distance of parameters cannot capture dynamics' dissimilarity

# Summary

$$\underset{\mathbf{f}_1, \dots, \mathbf{f}_n}{\text{minimize}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n L_i(\mathbf{f}_i; \{\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{\tau_i}^{(i)}\})}_{\text{original objective function to learn } \mathbf{f}_1, \dots, \mathbf{f}_n} + \lambda \cdot \underbrace{\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=i+1}^n \frac{d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j)}{d_{\mathbf{Z}}(Z_i, Z_j)}}_{\text{regularization term to make } d_{\text{DS}}(\mathbf{f}_i, \mathbf{f}_j) \propto d_{\mathbf{Z}}(Z_i, Z_j)}$$

e.g., squared loss  $L_i = \sum_{t=1}^{\tau_i-1} \|\mathbf{f}_i(\mathbf{x}_t^{(i)}) - \mathbf{x}_{t+1}^{(i)}\|_2^2$

- **Multi-task learning** formulation for dynamical systems
  - to utilize **side information**
- We adopt the **operator-theoretic metric** [Ishikawa+ 18] for measuring dissimilarity between dynamical systems,
  - which can be estimated from trajectories during training
  - good for utilizing side information

