

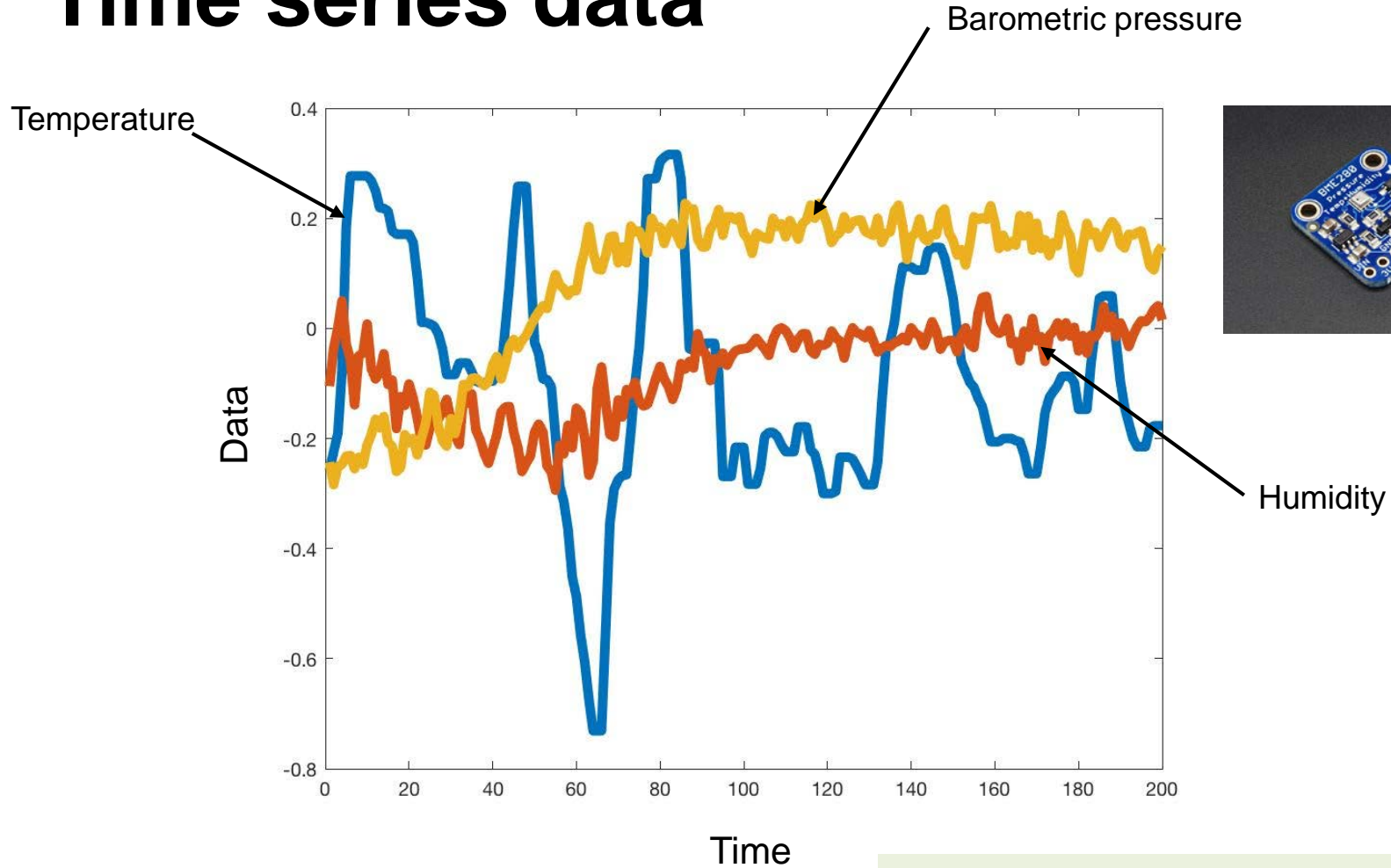
Modeling time series with hidden Markov models

Advanced Machine learning
2017

Nadia Figueroa, Jose Medina and Aude Billard



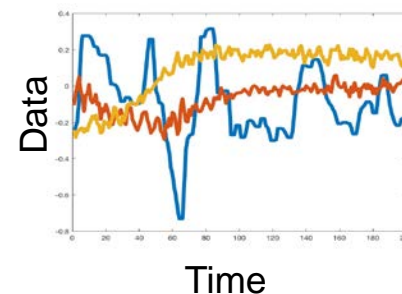
Time series data



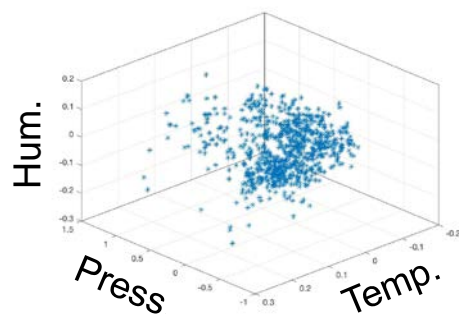
What's going on here?

Time series data

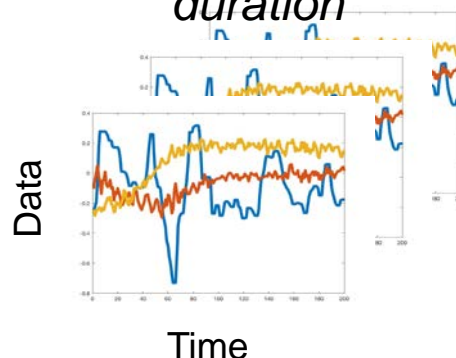
What's the problem setting?



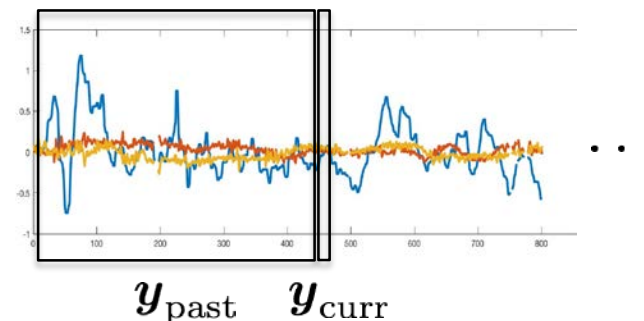
We don't care about time ...



We have several trajectories with identical duration



We have unstructured trajectory(ies)!



$$y = f(x)$$

Explicit time dependency

$$y = f(x, t)$$

$$y = f(t)$$

Consider dependency on the past

$$y_{\text{curr}} = f(y_{\text{past}})$$

Too complex!



LASA

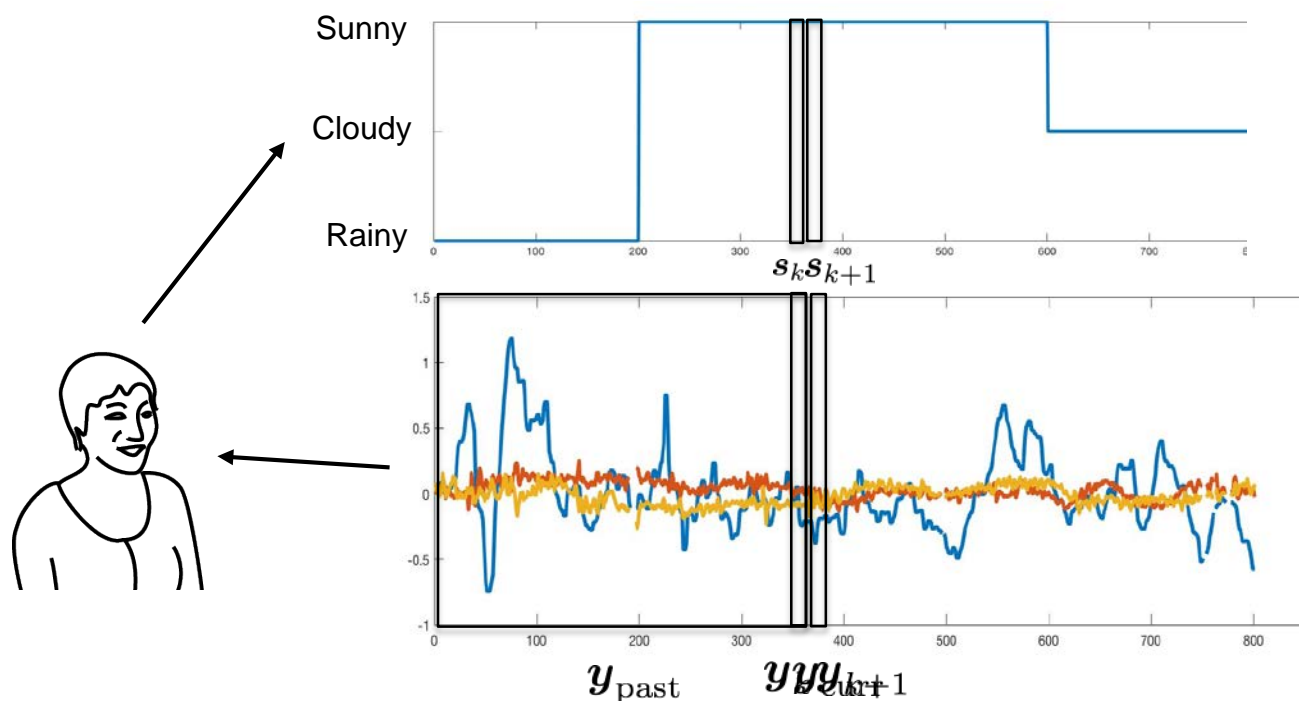
Learning Algorithms and
Systems Laboratory



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Unstructured time series data

How to simplify this problem?



$$s_{k+1} = f_{\text{dyn}}(s_k)$$

Markov assumption

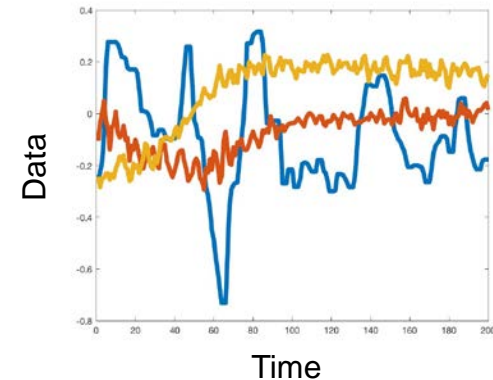
Consider dependency on the past

$$y_{k+1} = f(y_k, y_{k-1}, \dots, y_{\text{past}})$$

Outline

First part (10:15 – 11:00):

- Recap on Markov chains
- Hidden Markov Model (HMM)
 - Recognition of time series
 - ML Parameter estimation



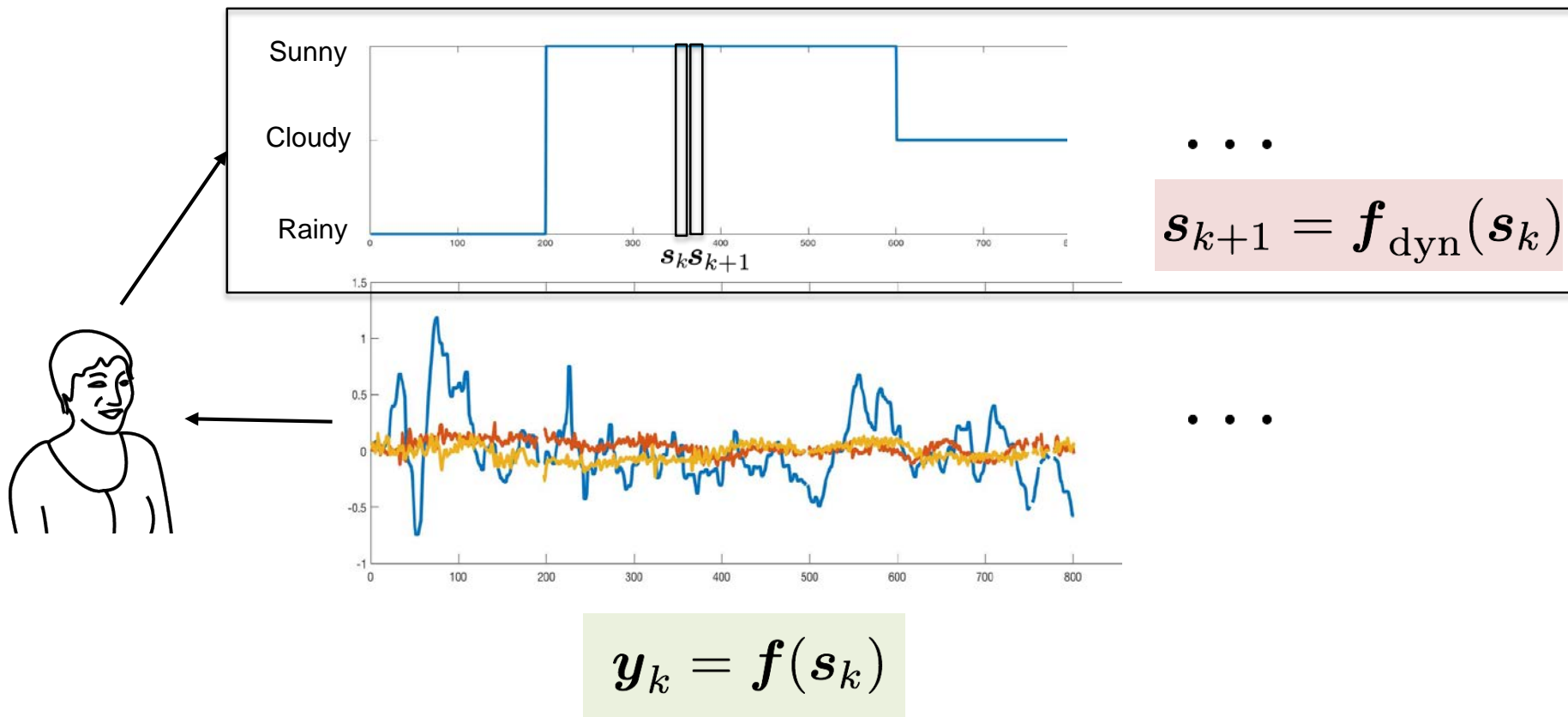
Second part (11:15 – 12:00):

- Time series segmentation
- Bayesian nonparametrics for HMMs

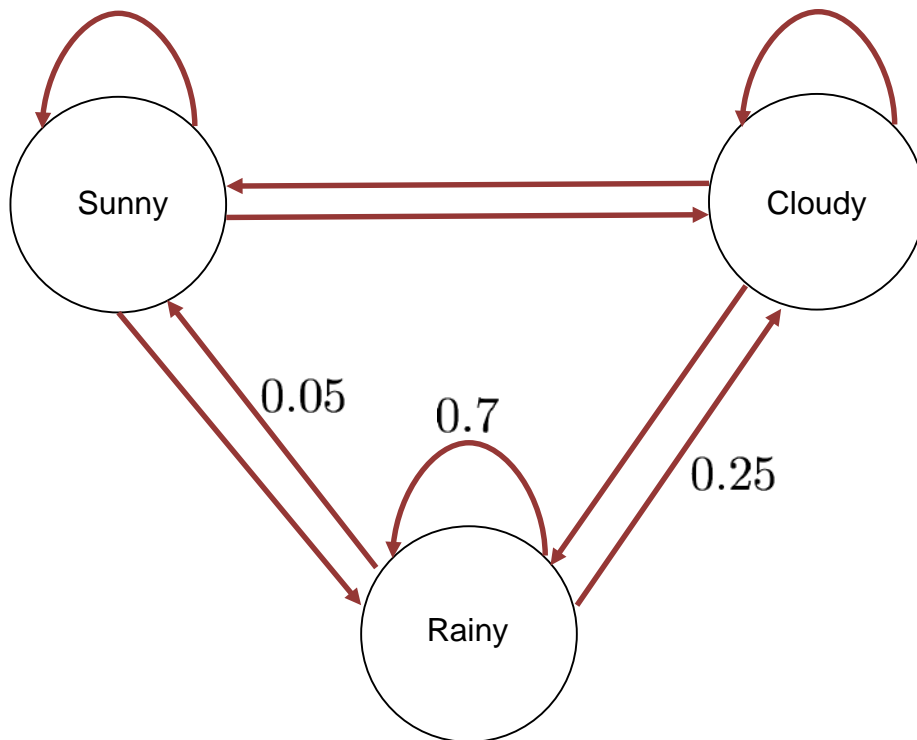
https://github.com/epfl-lasa/ML_toolbox



Outline first part



Markov chains



$$P(s_{k+1} = \text{Sunny} | s_k = \text{Rainy}) = 0.05$$

$$P(s_{k+1} = \text{Rainy} | s_k = \text{Rainy}) = 0.7$$

$$P(s_{k+1} = \text{Cloudy} | s_k = \text{Rainy}) = 0.25$$

$$s_{k+1} = f_{\text{dyn}}(s_k)$$

$$\lambda =$$

Transition

matrix

	Sunny	Cloudy	Rainy
Sunny	0.9	0.08	0.02
Cloudy	0.25	0.7	0.05
Rainy	0.05	0.25	0.7

Initial probabilities

	Sunny	Cloudy	Rainy
	0.33	0.33	0.33



Likelihood of a Markov chain



$O = \text{Sunny} \rightarrow \text{Sunny} \rightarrow \text{Cloudy}$

$\lambda =$

$$P(O|\lambda)?$$

$$\begin{aligned} P(s_0 = \text{Sunny}, s_1 = \text{Sunny}, s_2 = \text{Cloudy}) &= \\ P(s_0 = \text{Sunny})P(s_1 = \text{Sunny}|s_0 = \text{Sunny}) &= \\ P(s_2 = \text{Cloudy}|s_1 = \text{Sunny}) &= 0.33 \cdot 0.9 \cdot 0.08 \end{aligned}$$

Transition

matrix

	Sunny	Cloudy	Rainy
Sunny	0.9	0.08	0.02
Cloudy	0.25	0.7	0.05
Rainy	0.05	0.25	0.7

Initial probabilities

	Sunny	Cloudy	Rainy
	0.33	0.33	0.33

$$P(O|\lambda) = P(s_0) \prod_{k=1}^T P(s_k | s_{k-1})$$



Learning Markov chains

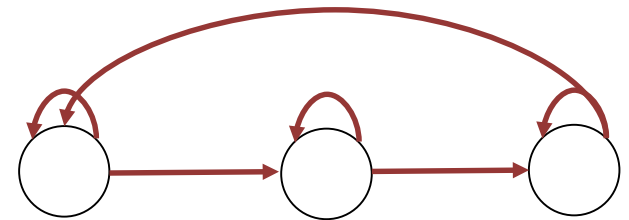


$O = \text{Sunny} \rightarrow \text{Sunny} \rightarrow \text{Cloudy}$

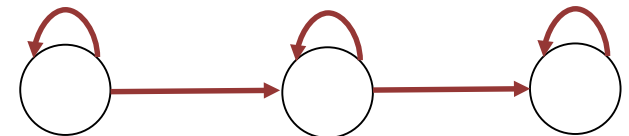
$$\max_{\lambda} \log P(O|\lambda)$$

Topologies

Periodic

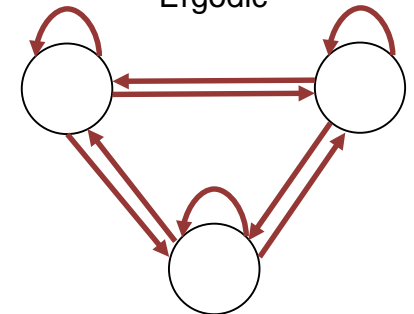


Left-to-right

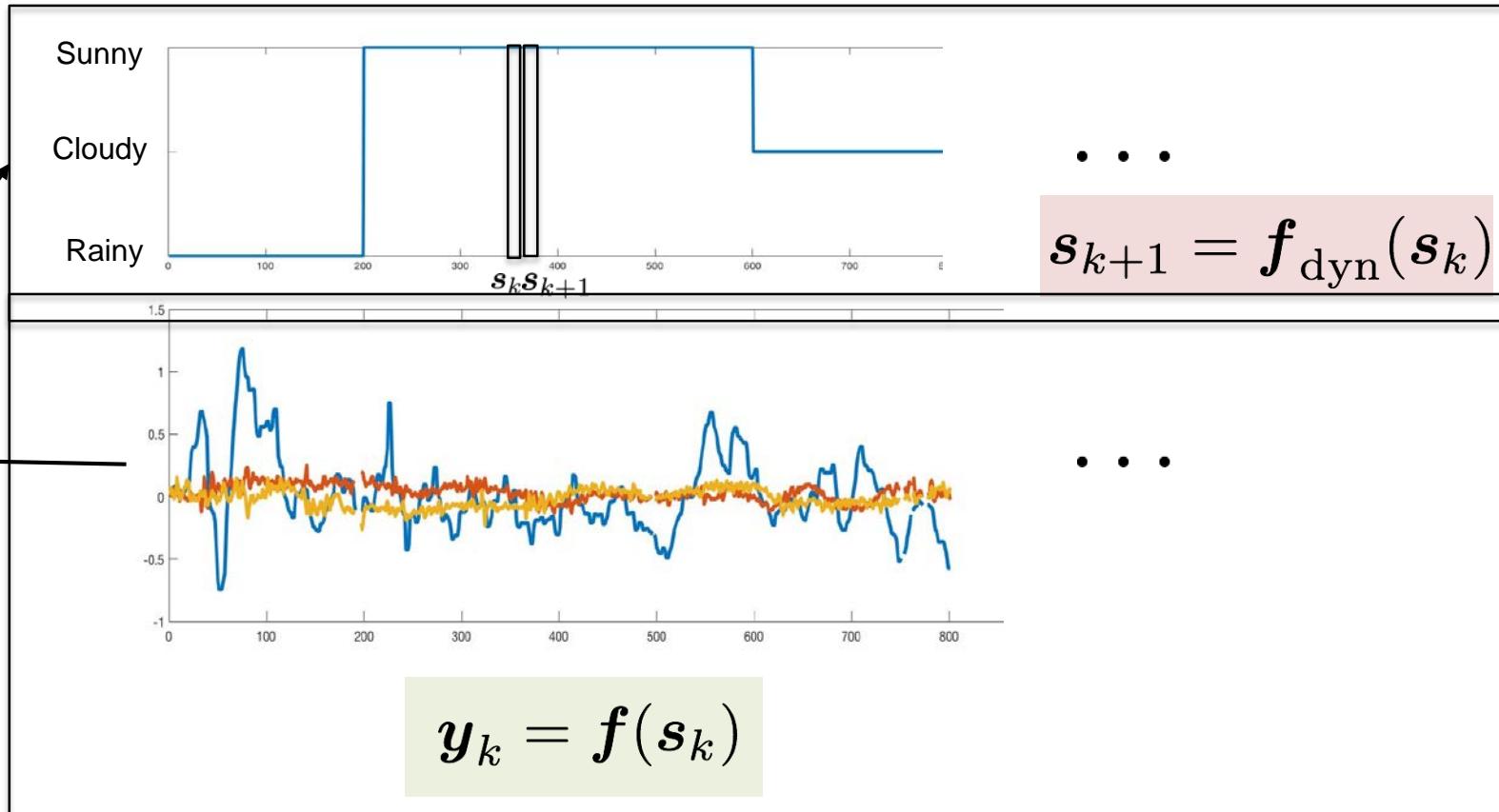


$$\hat{P}(s_{k+1} = \text{Sunny} | s_k = \text{Rainy}) = \frac{\text{number of times we transition Rainy} \rightarrow \text{Sunny}}{\text{number of times we observe Sunny}}$$

Ergodic

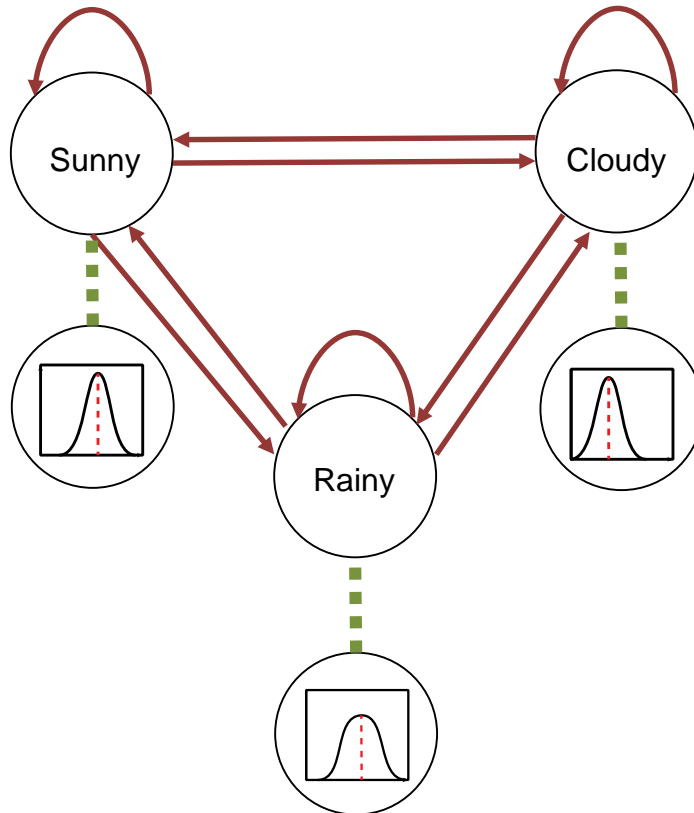


Outline first part



Hidden Markov model

$$s_{k+1} = f_{\text{dyn}}(s_k)$$



Transition

	Sunny	Cloudy	Rainy
Sunny	0.9	0.08	0.02
Cloudy	0.25	0.7	0.05
Rainy	0.05	0.25	0.7

Initial probabilities

Sunny	Cloudy	Rainy
0.33	0.33	0.33

$$y_k = f(s_k)$$

$$P(y_k | s_k = \text{Sunny}) = \mathcal{N}(y_k; \mu_{\text{Sunny}}, \Sigma_{\text{Sunny}})$$

$$P(y_k | s_k = \text{Cloudy}) = \mathcal{N}(y_k; \mu_{\text{Cloudy}}, \Sigma_{\text{Cloudy}})$$

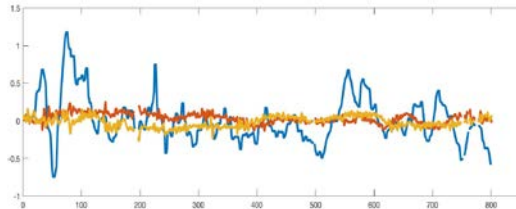
$$P(y_k | s_k = \text{Rainy}) = \mathcal{N}(y_k; \mu_{\text{Rainy}}, \Sigma_{\text{Rainy}})$$

Likelihood of an HMM

$$s_{k+1} = f_{\text{dyn}}(s_k)$$



$O =$



$y_1 \cdots y_T$

Transition

	Sunny	matrix	Rainy
Sunny	0.9	0.08	0.02
Cloudy	0.25	0.7	0.05
Rainy	0.05	0.25	0.7

Initial probabilities

Sunny	Cloudy	Rainy
0.33	0.33	0.33

$$P(O|\lambda) = P(y_1 \cdots y_T|\lambda) =$$

$$\sum_{s_1 \dots s_T \in \mathbb{D}} P(y_1 \cdots y_T, s_1 \cdots s_T|\lambda)$$

$\mathbb{D} = \{\text{Sunny, Cloudy, Rainy}\}$

Forward variable

$$\alpha(s_k) = P(y_1 \cdots y_k, s_k)$$

$$P(O|\lambda) = \sum_{s_T} \alpha(s_T)$$

$$y_k = f(s_k)$$

$$P(y_k | s_k = \text{Sunny}) = \mathcal{N}(y_k; \mu_{\text{Sunny}}, \Sigma_{\text{Sunny}})$$

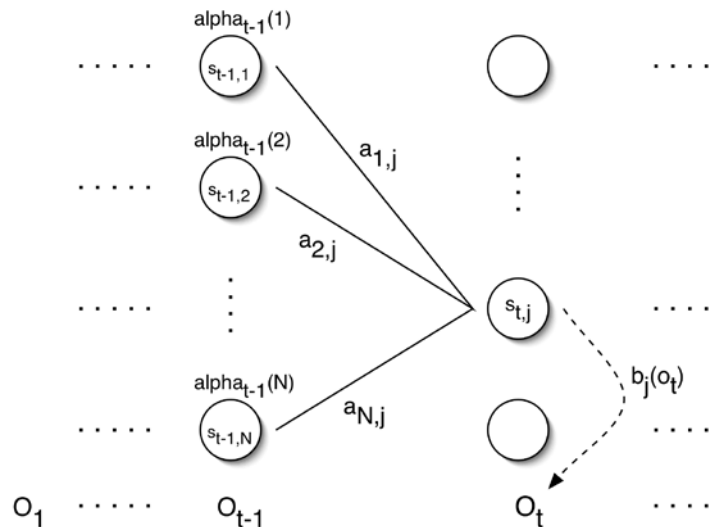
$$P(y_k | s_k = \text{Cloudy}) = \mathcal{N}(y_k; \mu_{\text{Cloudy}}, \Sigma_{\text{Cloudy}})$$

$$P(y_k | s_k = \text{Rainy}) = \mathcal{N}(y_k; \mu_{\text{Rainy}}, \Sigma_{\text{Rainy}})$$



Likelihood of an HMM

$$\alpha(s_k) = \sum_{s_{k-1}} \alpha(s_{k-1}) P(y_k | s_k) P(s_k | s_{k-1})$$



Forward variable

$$\alpha(s_k) = P(y_1 \cdots y_k, s_k)$$

$$P(O|\lambda) = \sum_{s_T} \alpha(s_T)$$

$$s_{k+1} = f_{\text{dyn}}(s_k)$$

Transition

	Sunny	matrix	Rainy
Sunny	0.9	0.08	0.02
Cloudy	0.25	0.7	0.05
Rainy	0.05	0.25	0.7

Initial probabilities

	Sunny	Cloudy	Rainy
	0.33	0.33	0.33

$$y_k = f(s_k)$$

$$P(y_k | s_k = \text{Sunny}) = \mathcal{N}(y_k; \mu_{\text{Sunny}}, \Sigma_{\text{Sunny}})$$

$$P(y_k | s_k = \text{Cloudy}) = \mathcal{N}(y_k; \mu_{\text{Cloudy}}, \Sigma_{\text{Cloudy}})$$

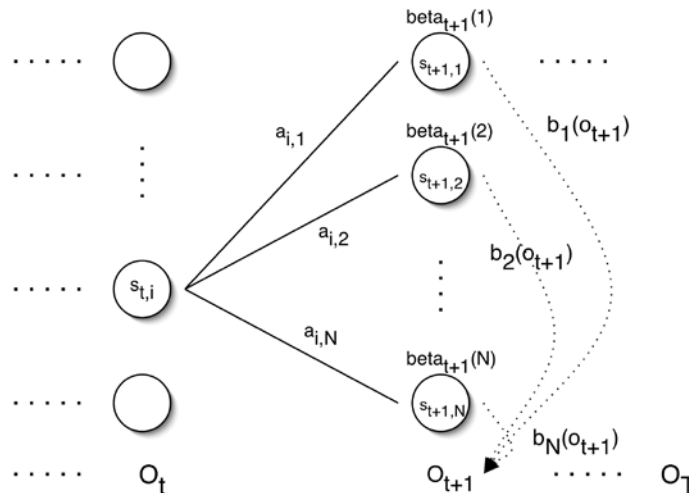
$$P(y_k | s_k = \text{Rainy}) = \mathcal{N}(y_k; \mu_{\text{Rainy}}, \Sigma_{\text{Rainy}})$$



Likelihood of an HMM

$$s_{k+1} = f_{\text{dyn}}(s_k)$$

$$\beta(s_k) = \sum_{s_{k+1}} \beta(s_{k+1}) P(\mathbf{y}_{k+1} | s_{k+1}) P(s_{k+1} | s_k)$$



Transition

	Sunny	matrix	Rainy
Sunny	0.9	0.08	0.02
Cloudy	0.25	0.7	0.05
Rainy	0.05	0.25	0.7

Initial probabilities

	Sunny	Cloudy	Rainy
	0.33	0.33	0.33

Backward variable

$$\beta(s_k) = P(\mathbf{y}_{k+1} \cdots \mathbf{y}_T | s_k)$$

$$P(O | \lambda) = \sum_{s_1} \beta(s_1) P(s_1)$$

$$\mathbf{y}_k = f(s_k)$$

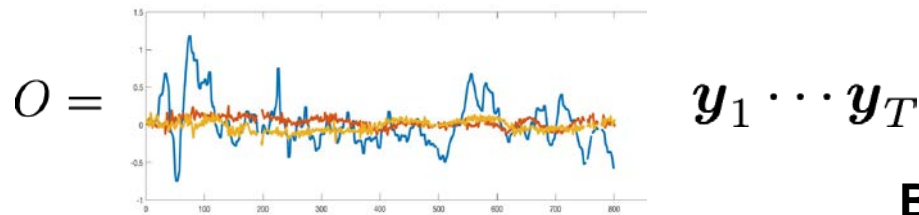
$$P(\mathbf{y}_k | s_k = \text{Sunny}) = \mathcal{N}(\mathbf{y}_k; \boldsymbol{\mu}_{\text{Sunny}}, \boldsymbol{\Sigma}_{\text{Sunny}})$$

$$P(\mathbf{y}_k | s_k = \text{Cloudy}) = \mathcal{N}(\mathbf{y}_k; \boldsymbol{\mu}_{\text{Cloudy}}, \boldsymbol{\Sigma}_{\text{Cloudy}})$$

$$P(\mathbf{y}_k | s_k = \text{Rainy}) = \mathcal{N}(\mathbf{y}_k; \boldsymbol{\mu}_{\text{Rainy}}, \boldsymbol{\Sigma}_{\text{Rainy}})$$



Learning an HMM



$$\max_{\lambda} \log P(O|\lambda)$$

Baum-Welch algorithm

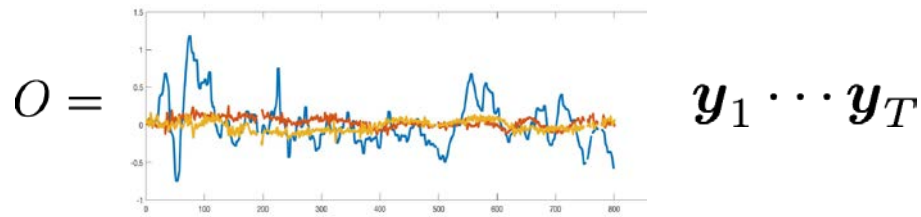
(Expectation-Maximization for HMMs)

- Iterative solution
- Converges to local minimum

Starting from an initial λ find a λ' such that $P(O|\lambda') \geq P(O|\lambda)$

- **E-step:** Given an observation sequence and a model, find the probabilities of the states to have produced those observations.
- **M-step:** Given the output of the E-step, update the model parameters to better fit the observations.

Learning an HMM



E-step:

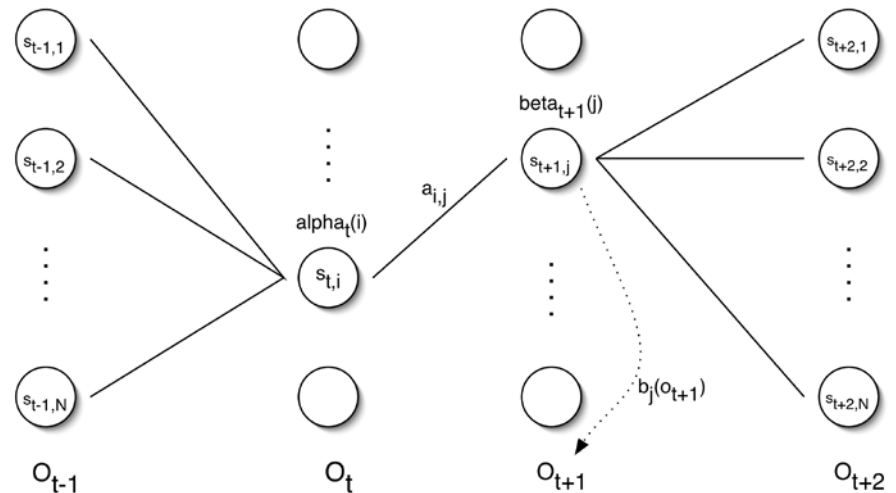
$$\xi_k(i, j) = \frac{\alpha(s_k = i)P(s_k = i, s_{k+1} = j)P(y_{k+1}|s_{k+1} = j)\beta(s_{k+1} = j)}{\sum_{s_k} \sum_{s_{k+1}} \alpha(s_k)P(s_k, s_{k+1})P(y_{k+1}|s_{k+1})\beta(s_{k+1})}$$

Probability of being in state i at time k and transition to state j

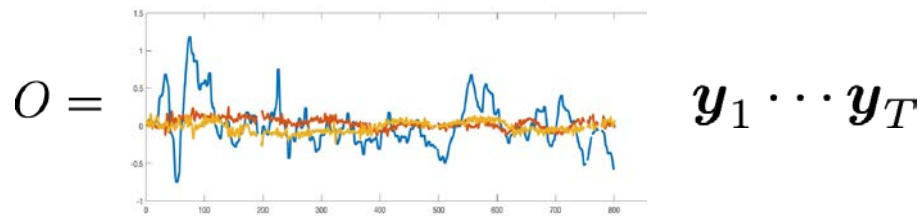
$$\xi_k(i, j) = P(s_k = i, s_k = j | O, \lambda)$$

Probability of being in state i at time k

$$\gamma_k(i) = \sum_{j \in s_{k+1}} \xi_k(i, j)$$



Learning an HMM



M-step:

$$\hat{P}(s_{k+1} = \text{Sunny} | s_k = \text{Rainy}) = \frac{\sum_{k=1}^T \xi_k(i, j)}{\sum_{k=1}^T \gamma_k(j)}$$

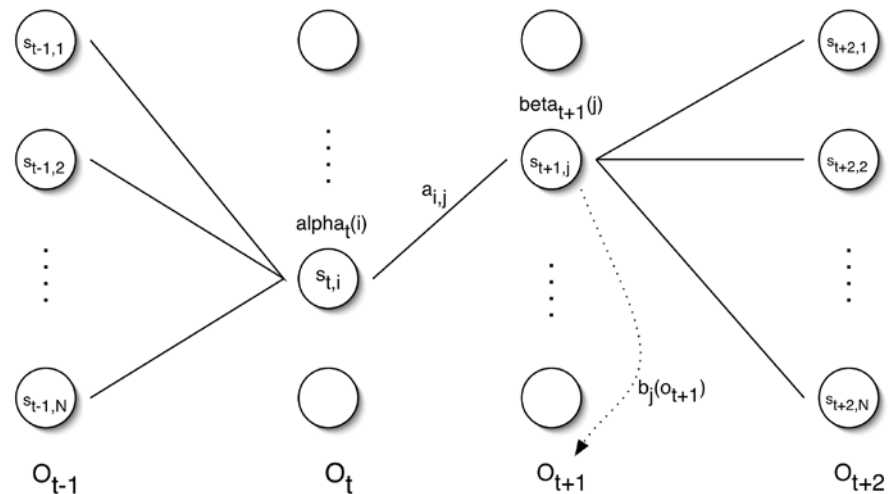
number of times we transition Rainy \rightarrow Sunny
number of times we observe Sunny

Probability of being in state i at time k and transition to state j

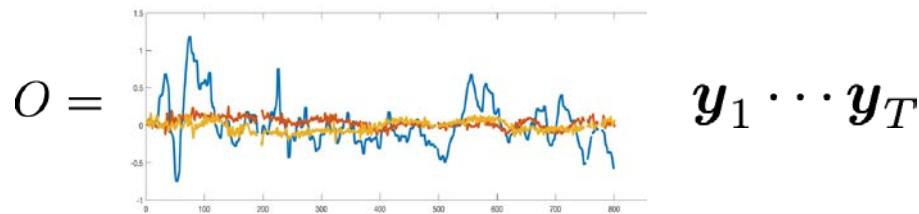
$$\xi_k(i, j) = P(s_k = i, s_k = j | O, \lambda)$$

Probability of being in state i at time k

$$\gamma_k(i) = \sum_{j \in s_{k+1}} \xi_k(i, j)$$



Learning an HMM



- HMM is a parametric technique (Fixed number of states, fixed topology)
- → Heuristics to determining the optimal number of states

X : dataset; N : number of datapoints; K : number of free parameters

- Akaike Information Criterion: $AIC = -2 \ln L + 2K$

- Bayesian Information Criterion: $BIC = -2 \ln L + K \ln(N)$

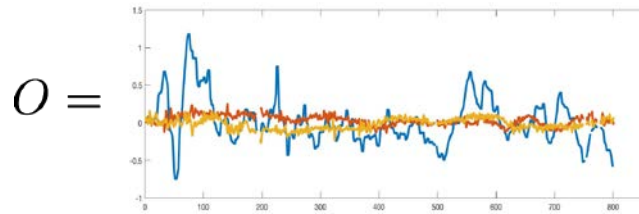
L : maximum likelihood of the model given K parameters

Choosing AIC versus BIC depends on the application:

→ Is the purpose of the analysis to make predictions, or to decide which model best represents reality?

AIC may have better predictive ability than BIC, but BIC finds a computationally more efficient solution.

Applications of HMMs



State estimation:

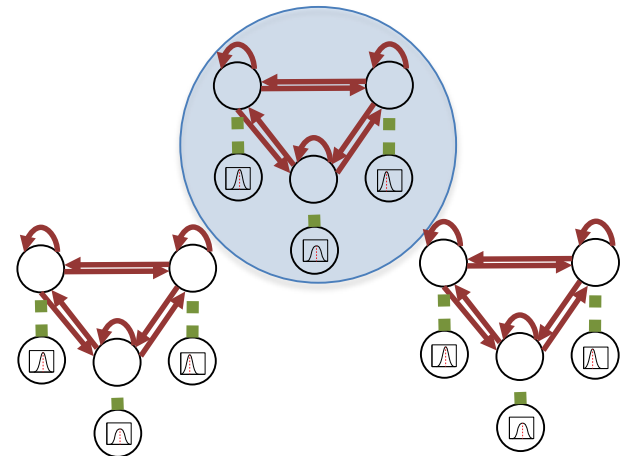
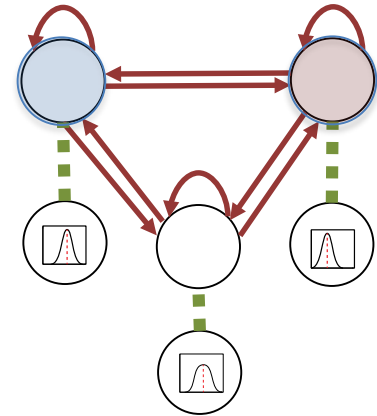
What is the most probable state/state sequence of the system?

Prediction:

What are the most probable next observations/state of the system?

Model selection:

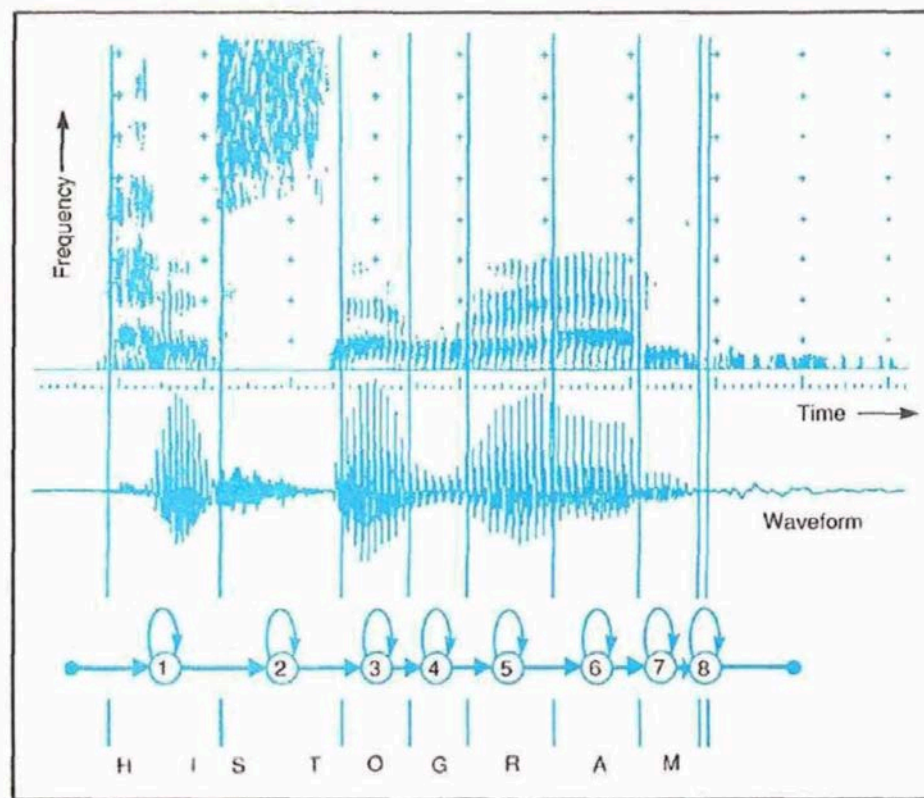
What is the most likely model that represents these observations?



Examples

Speech recognition:

- Left-to-right model
- States are phonemes
- Observations in frequency domain

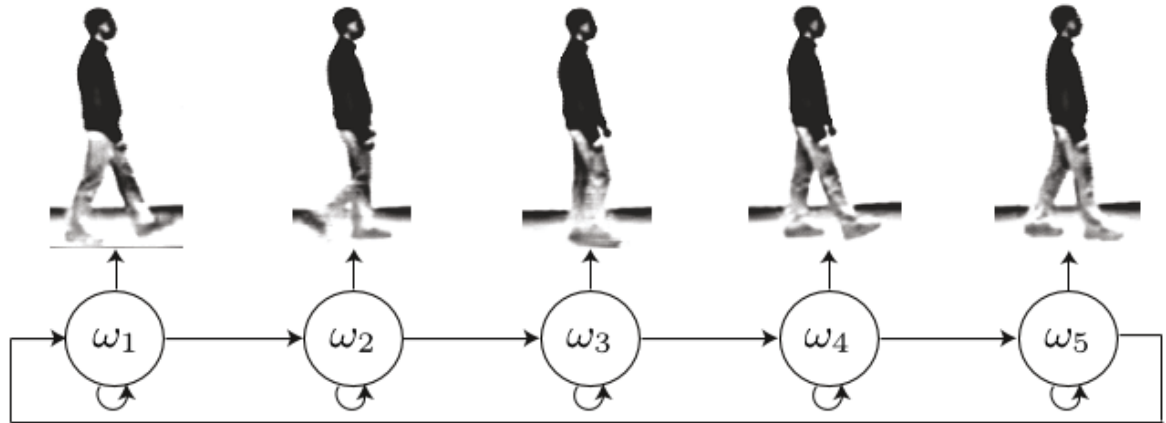


D.B. Paul., Speech Recognition Using Hidden Markov Models, The Lincoln laboratory journal, 1990

Examples

Motion prediction:

- Periodic model
- Observations are observed joints
- Simulate/predict walking patterns



Karg, Michelle, et al. "Human movement analysis: Extension of the f-statistic to time series using hmm." *Systems, Man, and Cybernetics (SMC), 2013 IEEE International Conference on*. IEEE, 2013.



Examples

Motion prediction:

- Left-to-right models
- Autonomous segmentation
- Recognition + prediction

An Incremental Learning and Prediction Framework and its Application in Human-Robot Cooperative Manipulation

J. R. Medina D. Lee S. Hirche

Examples

Motion prediction:

- Left-to-right models
- Autonomous segmentation
- Recognition + prediction

**An Experience-Driven Robotic Assistant
Acquiring Human Knowledge
to Improve Joint Manipulation**

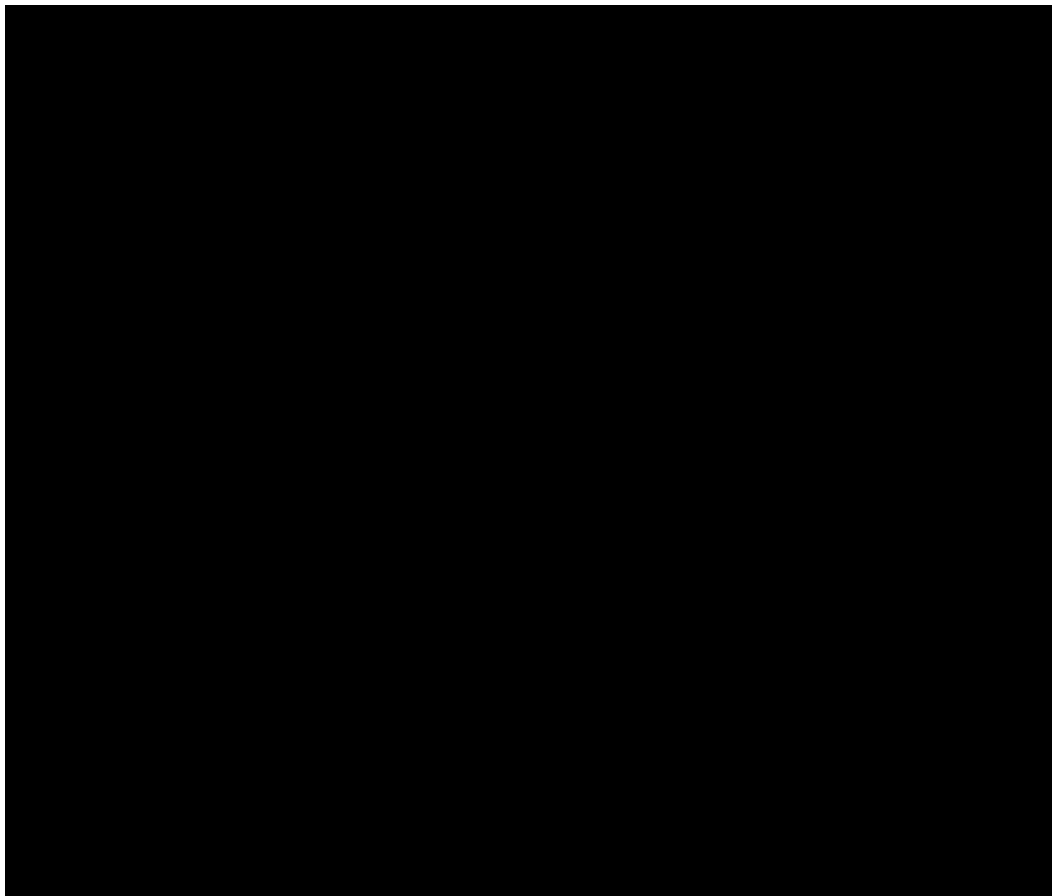
J.R.Medina M.Lawitzky A.Mörtl
D. Lee S.Hirche



Examples

Motion prediction:

- Left-to-right model
- Each state is a dynamical system



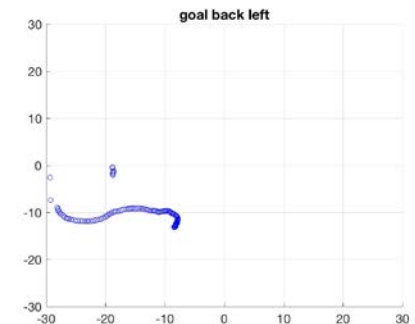
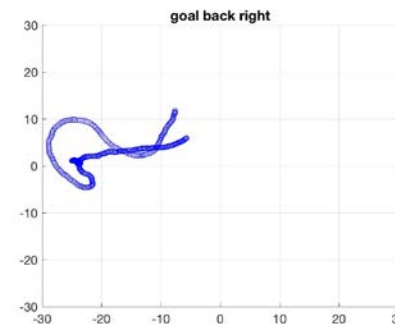
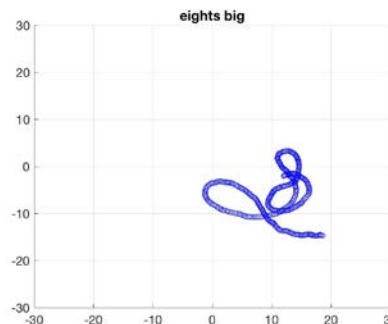
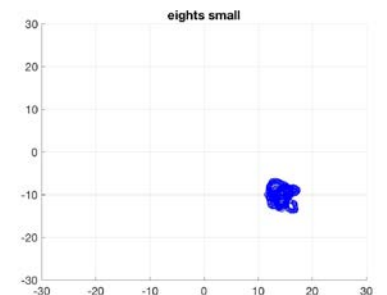
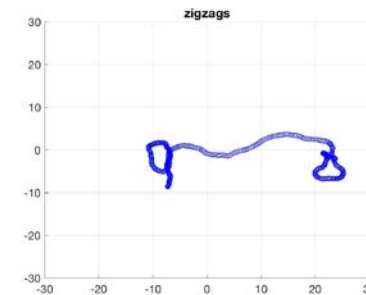
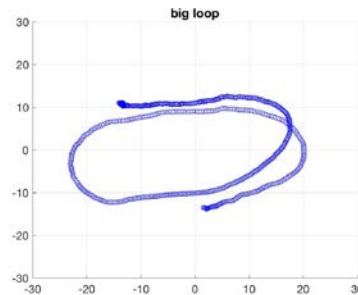
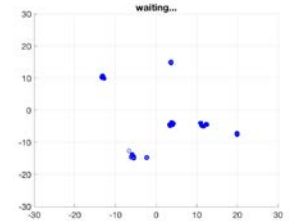
Examples

Motion recognition:

- Recognition of most likely motion and prediction of next step.

Toy training set

- 1 player
- 7 actions
- 1 Hidden Markov model per action



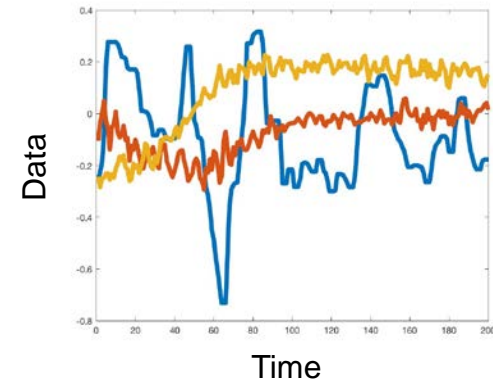
MATLAB demo



Outline

First part (10:15 – 11:00):

- Recap on Markov chains
- Hidden Markov Model (HMM)
 - Recognition of time series
 - ML Parameter estimation



Second part (11:15 – 12:00):

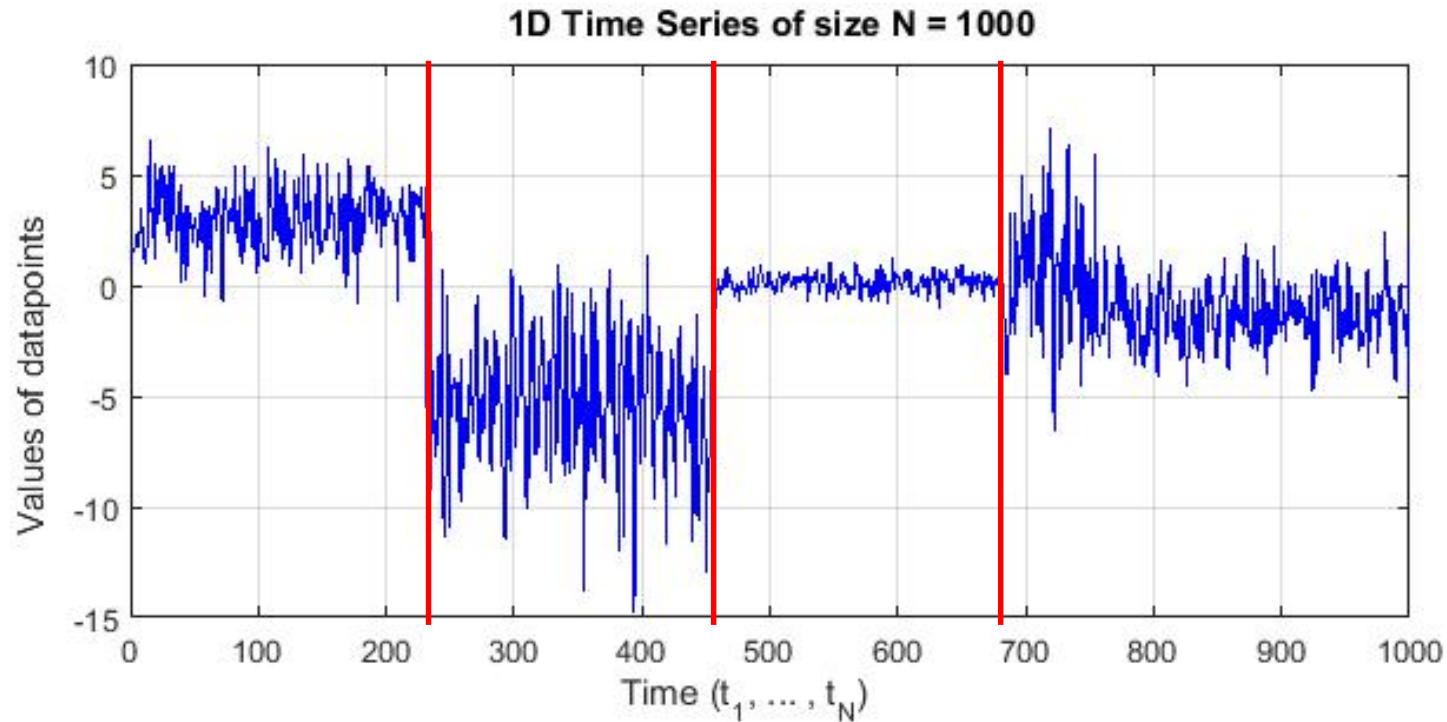
- Time series segmentation
- Bayesian non-parametrics for HMMs

https://github.com/epfl-lasa/ML_toolbox



Time series Segmentation

Times-series = Sequence of discrete segments



$O =$

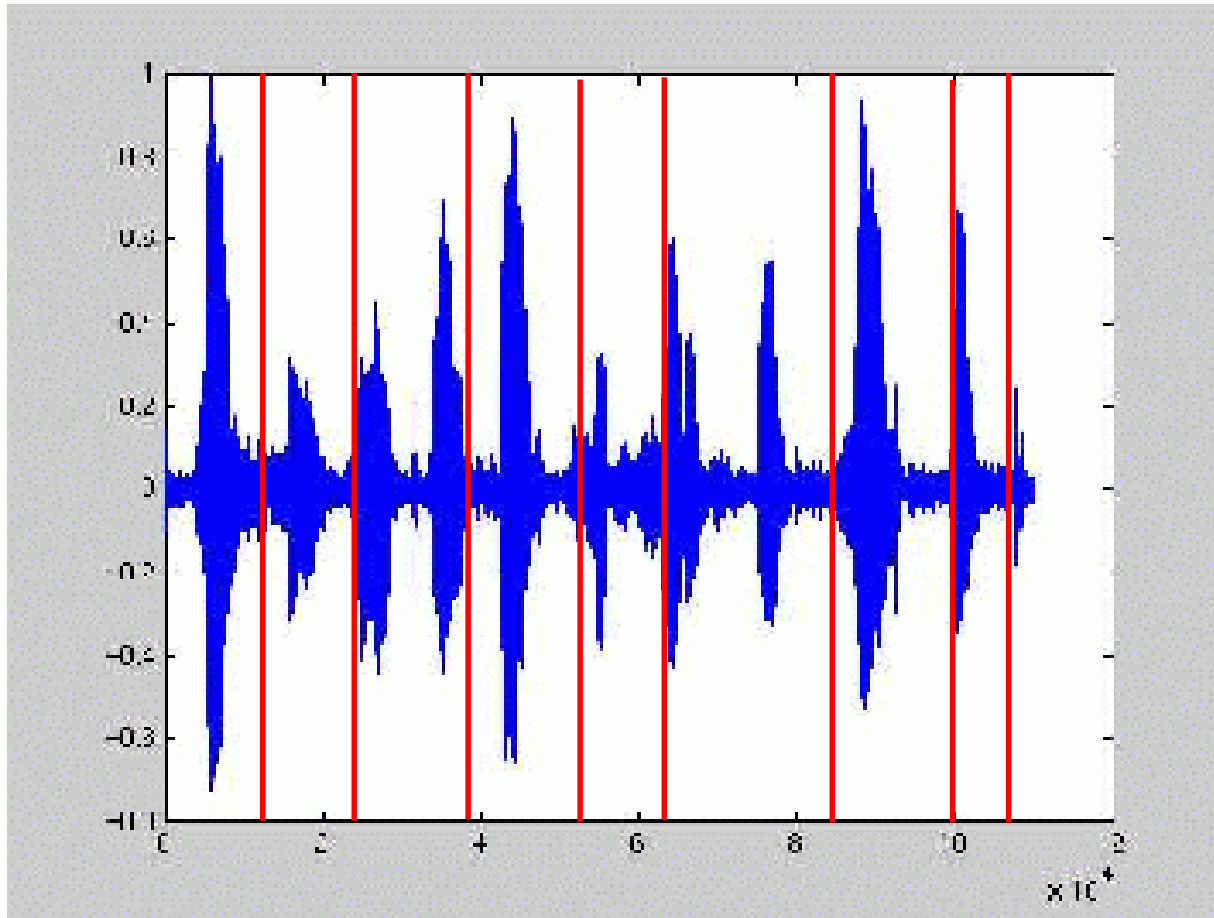
$y_1 \cdots y_T$

Why is this an important problem?

Segmentation of Speech Signals

Segmenting a continuous speech signal into sets of distinct words.

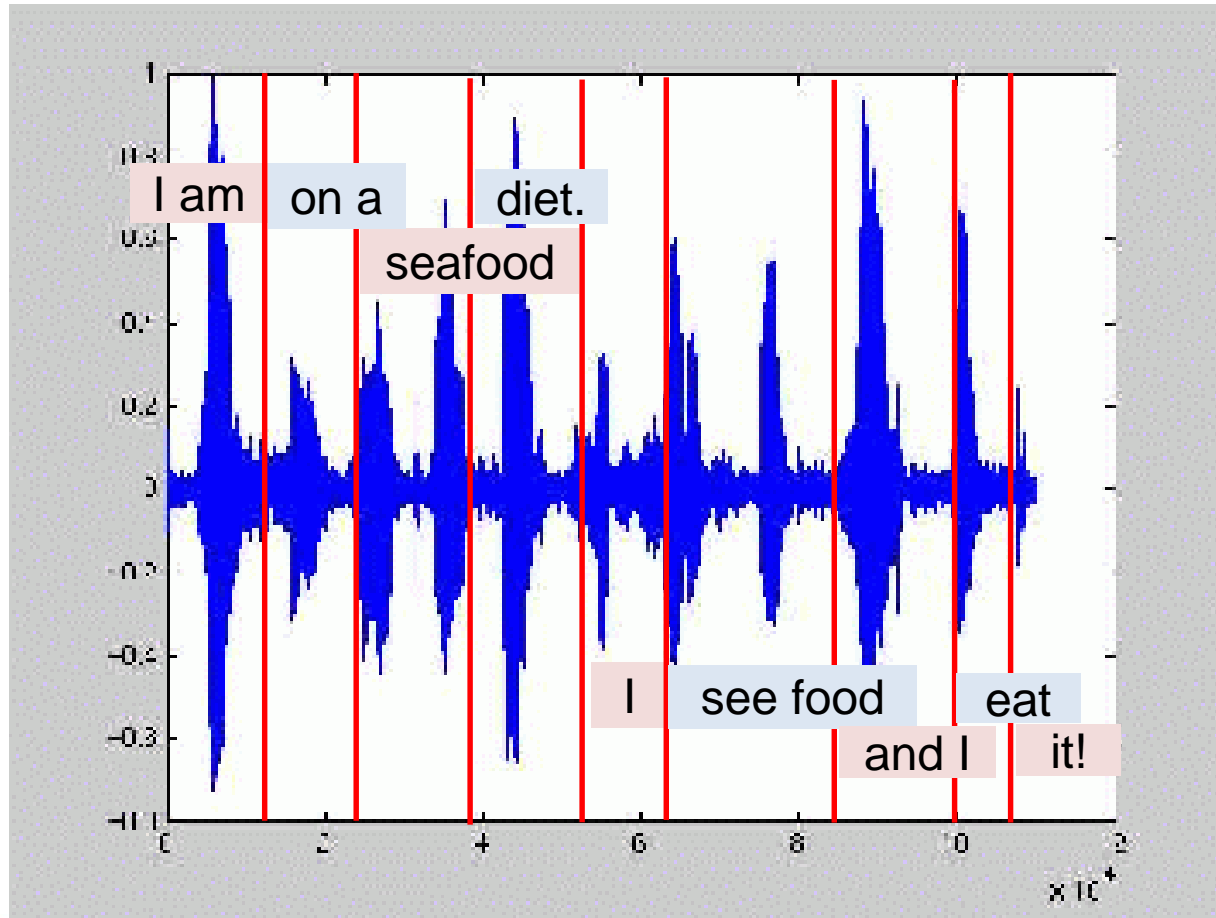
$O =$



$y_1 \cdots y_T$

Segmentation of Speech Signals

Segmenting a continuous speech signal into sets of distinct words.



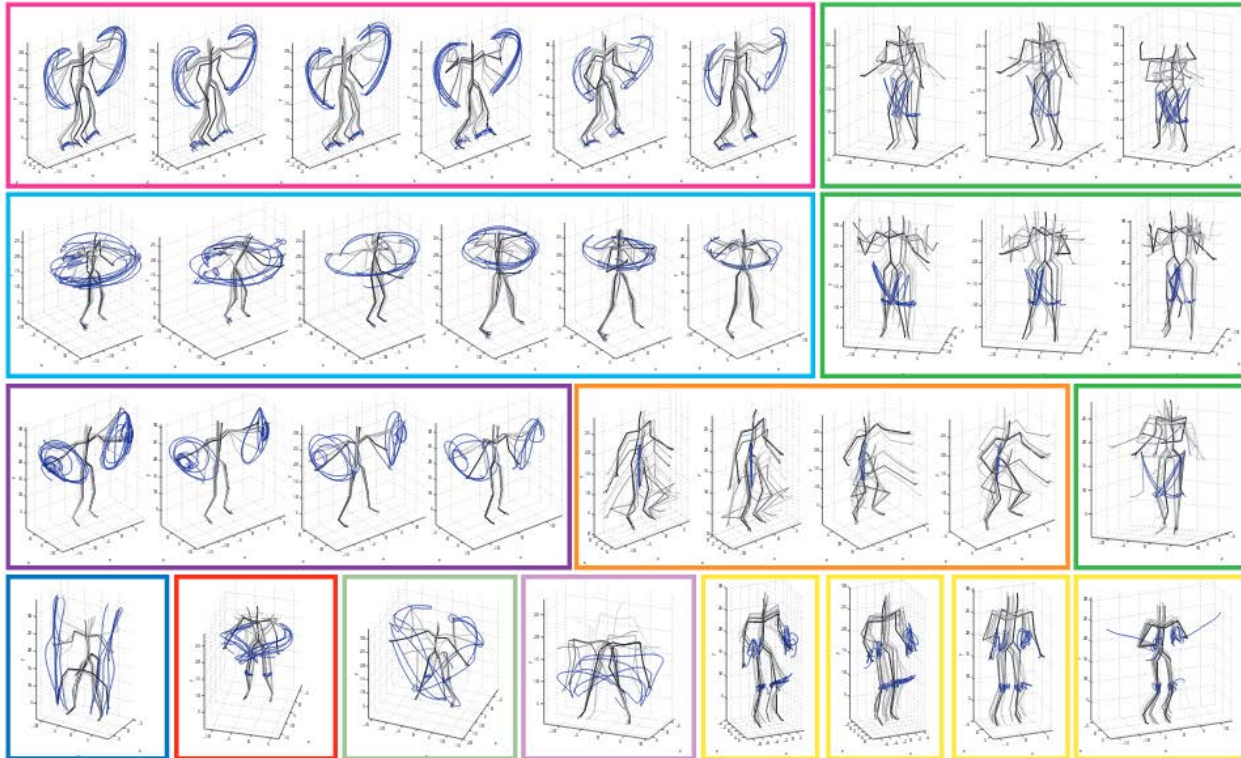
$O =$

$y_1 \cdots y_T$

Segmentation of Human Motion Data

Segmentation of Continuous Motion Capture data from exercise routines into motion categories

Jumping
Jacks



Knee
Raises

Arm
Circles

Squats

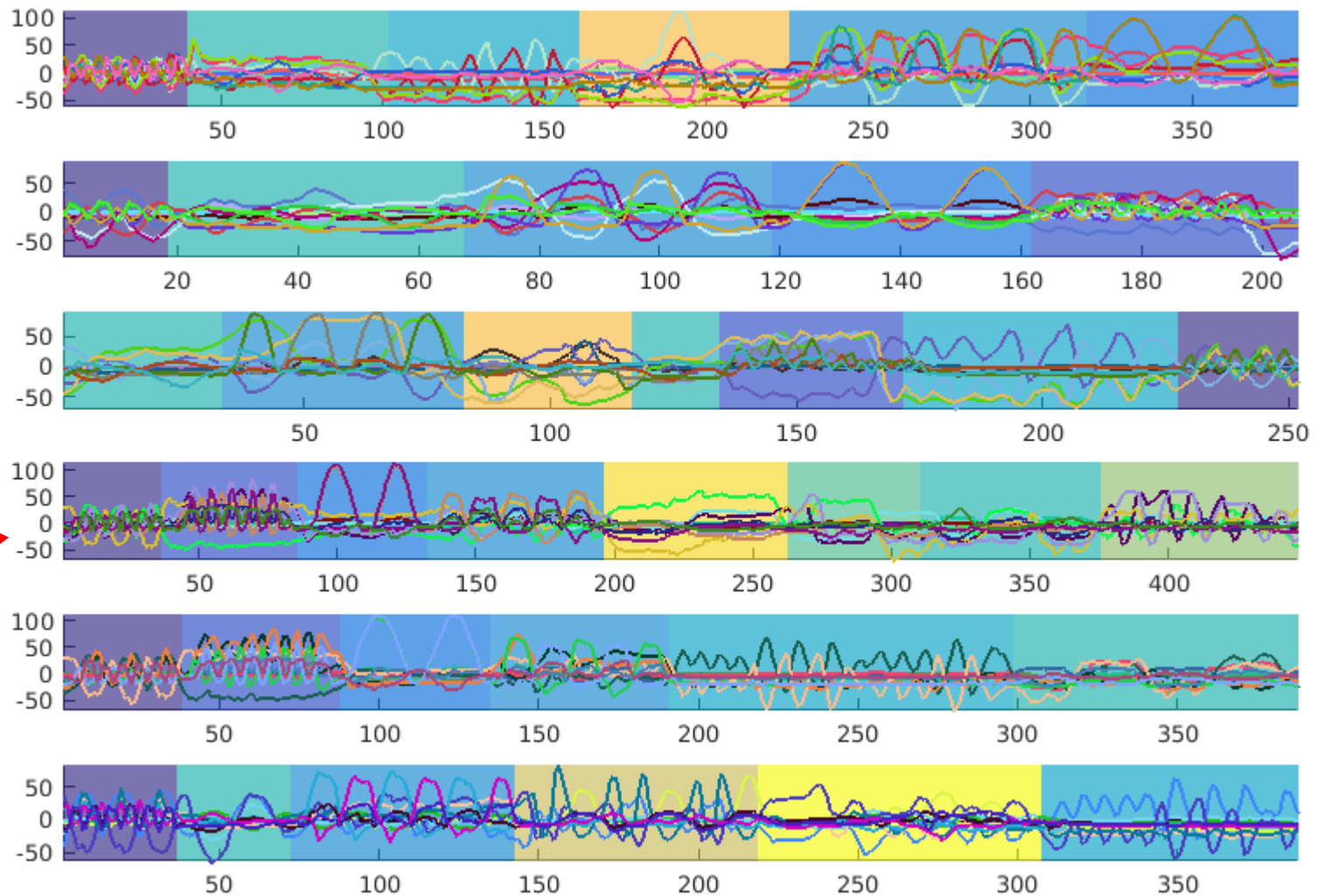
Emily Fox et al., Sharing Features among Dynamical Systems with Beta Processes, NIPS, 2009

Segmentation in Human Motion Data

12 Variables

- Torso position
- Waist Angles (2)
- Neck Angle
- Shoulder Angles
- ..

$$\mathbf{y}_i \in \mathbb{R}^{12}$$



Emily Fox et al., Sharing Features among Dynamical Systems with Beta Processes, NIPS, 2009



Segmentation in Robotics

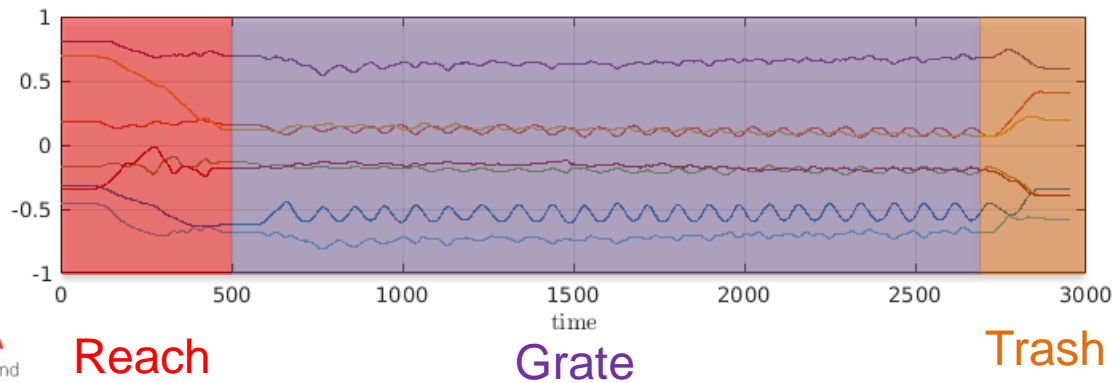
Learning Complex Sequential Tasks from Demonstration

7 Variables

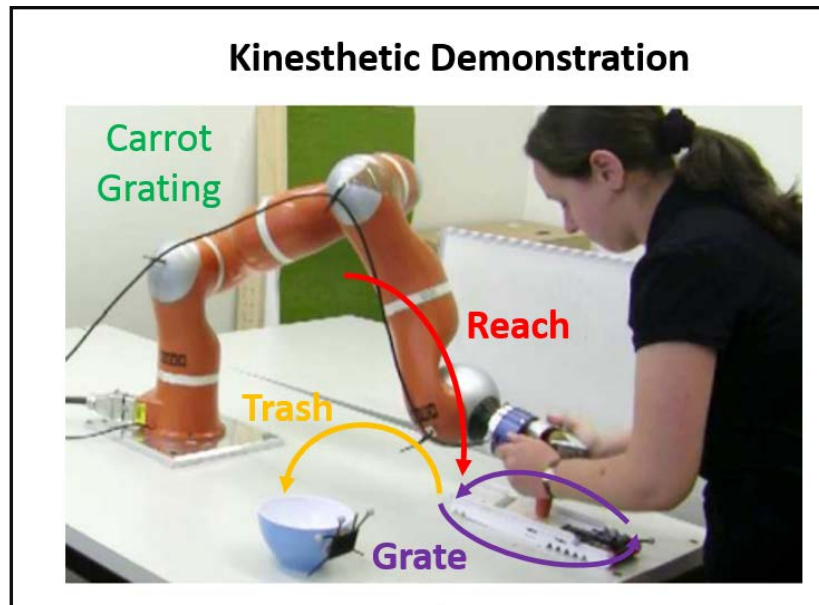
- Position
 x, y, z
- Orientation
 q_i, q_j, q_k, q_w

$$\mathbf{y}_i \in \mathbb{R}^7$$

$O =$



$\mathbf{y}_1 \cdots \mathbf{y}_T$



HMM for Time series Segmentation

Assumptions:

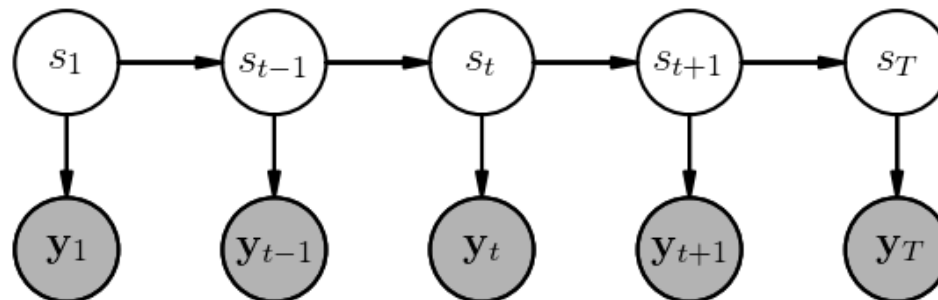
- The time-series has been generated by a system that **transitions** between a set of **hidden states**:

$$s \in \{1, \dots, K\}$$

- At each time step, a sample is drawn from an **emission model** associated to the current **hidden state**:

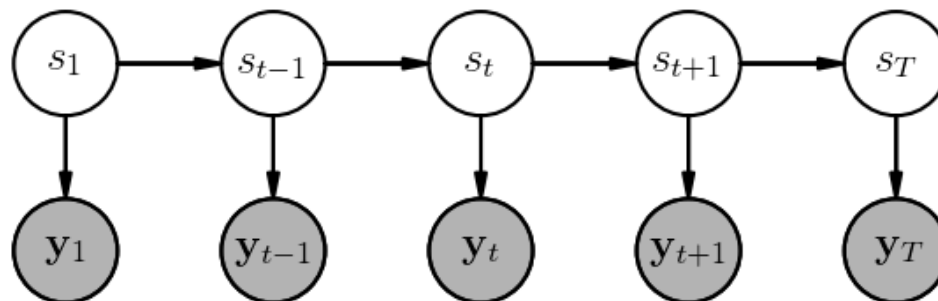
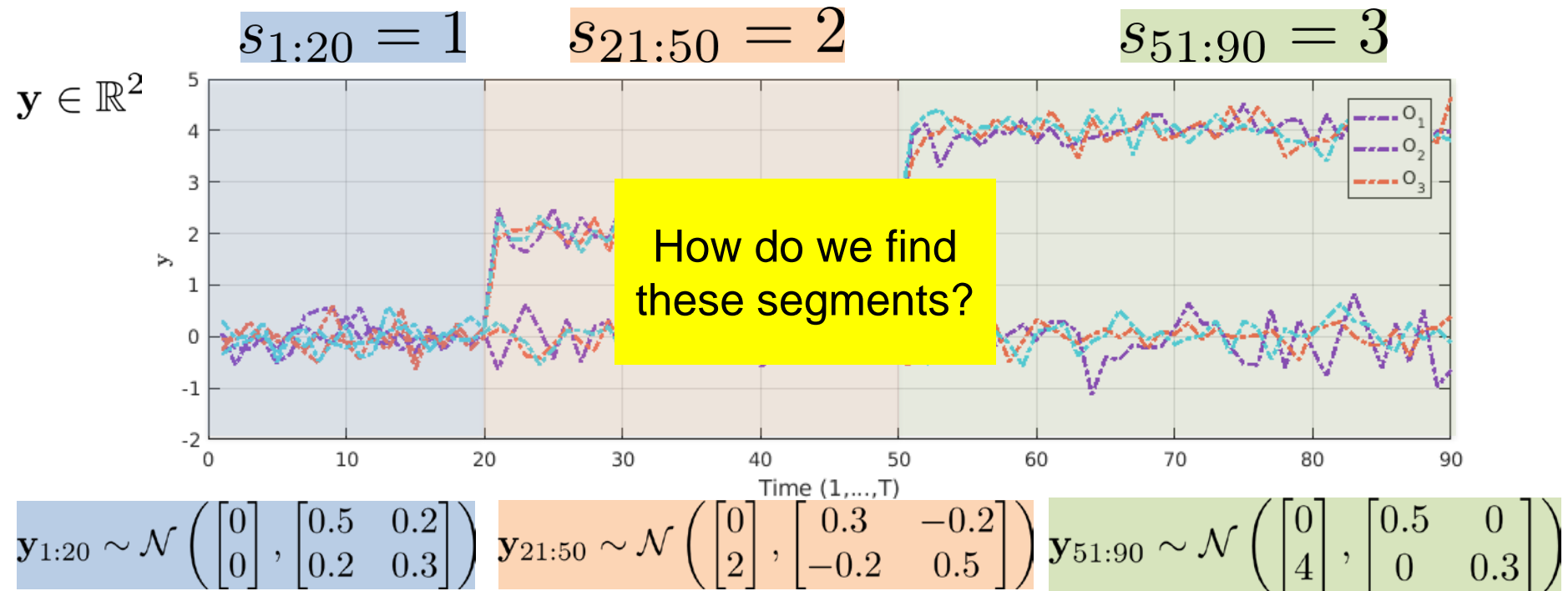
$$y_t | s_t = k \sim \mathcal{N}(\theta_{s_t})$$

$$\text{with } \theta_{s_t} = \{\mu_k, \Sigma_k\}$$





HMM for Time series Segmentation



HMM for Time series Segmentation

Steps for Segmentation with HMM:

1. Learn the HMM parameters through Maximum Likelihood Estimate (MLE):

$$\lambda = \{\pi, A, \Theta\}$$

Initial State Transition Emission Model
Probabilities Matrix Parameters

HMM Likelihood

$$P(O|\lambda) = P(\mathbf{y}_1 \cdots \mathbf{y}_T | \lambda) = \sum_{\mathbf{s}_1 \cdots \mathbf{s}_T \in \mathbb{D}} P(\mathbf{y}_1 \cdots \mathbf{y}_T, \mathbf{s}_1 \cdots \mathbf{s}_T | \lambda)$$

$$\max_{\lambda} \log P(O|\lambda)$$

Baum-Welch algorithm

(Expectation-Maximization for HMMs)

- Iterative solution
- Converges to local minimum

Hyper-parameter:
Number of states possible K



HMM for Time series Segmentation

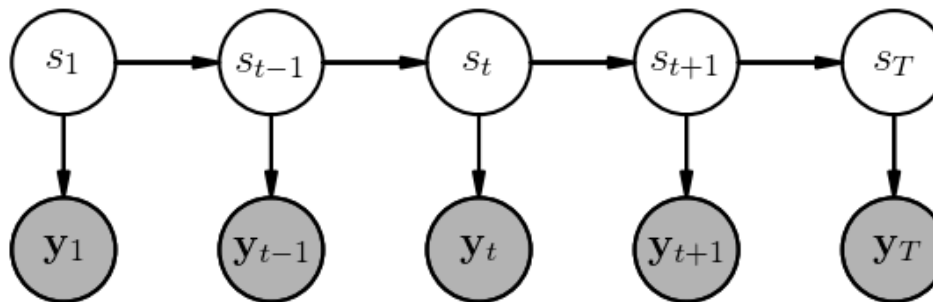
Steps for Segmentation with HMM:

2. Find the most probable sequence of states generating the observations through the **Viterbi algorithm**:

$$S^* = \arg \max_S p(Y, S | \Theta)$$

HMM Joint Probability Distribution

$$p(Y, S | \Theta) = p(s_1 | \pi) \left[\prod_{t=2}^T p(s_t | s_{t-1}, A) \right] \prod_{t=1}^T p(y_t | s_t, \Theta)$$



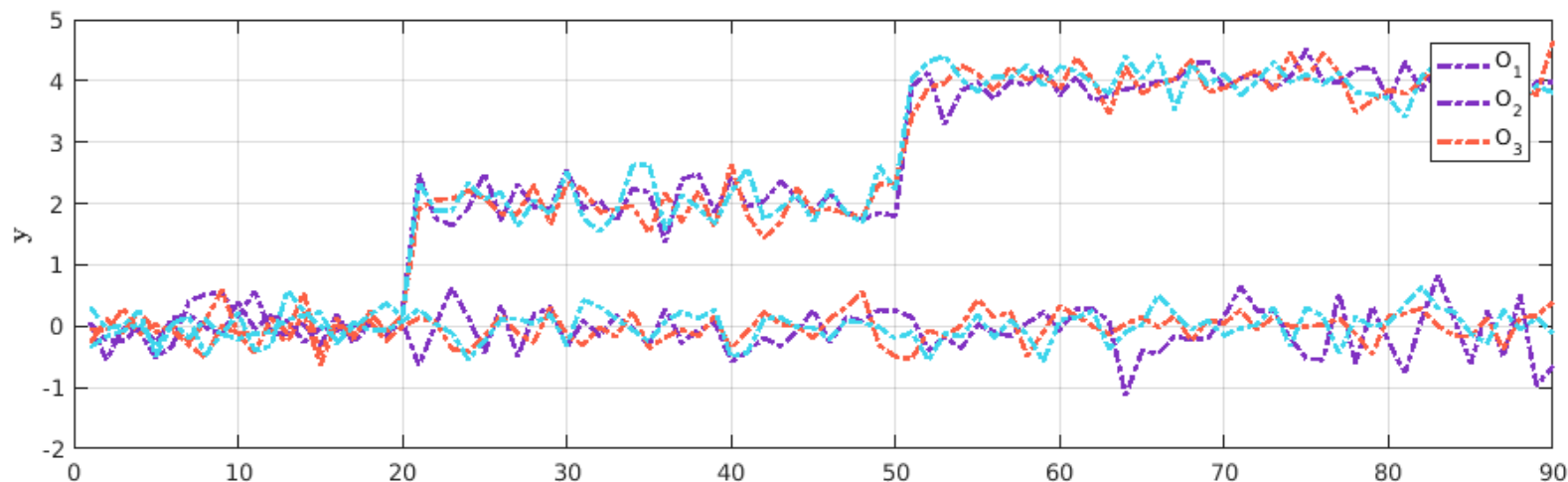
LASA

Learning Algorithms and
Systems Laboratory

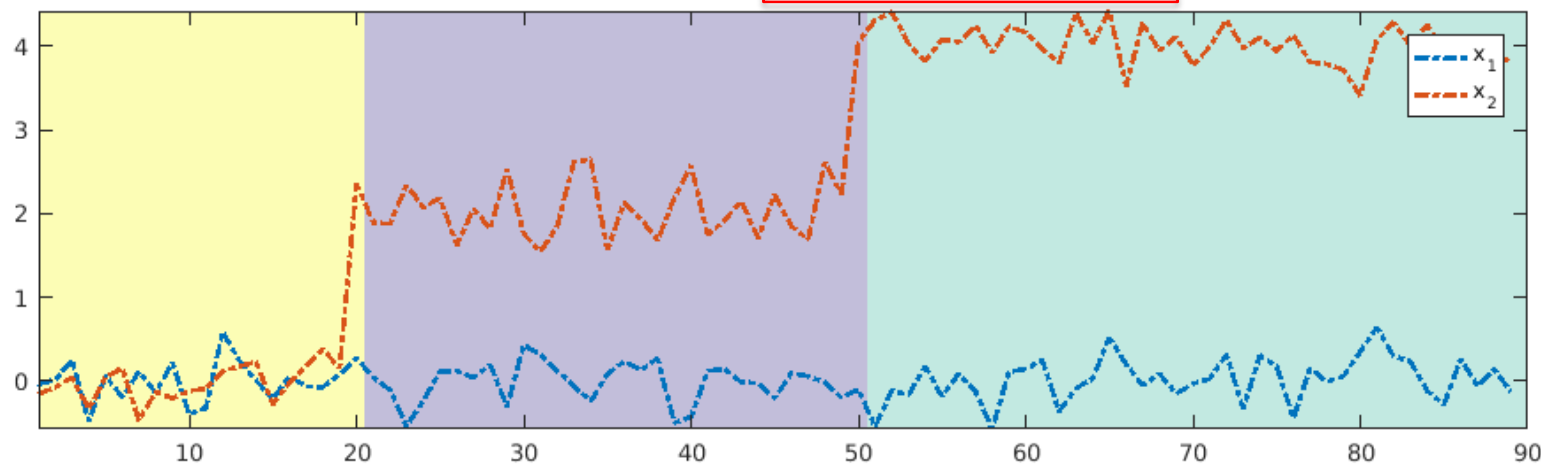


HMM for Time series Segmentation

$y \in \mathbb{R}^2$



Segmented Data, $K:3$, loglik:-52.0355



LASA

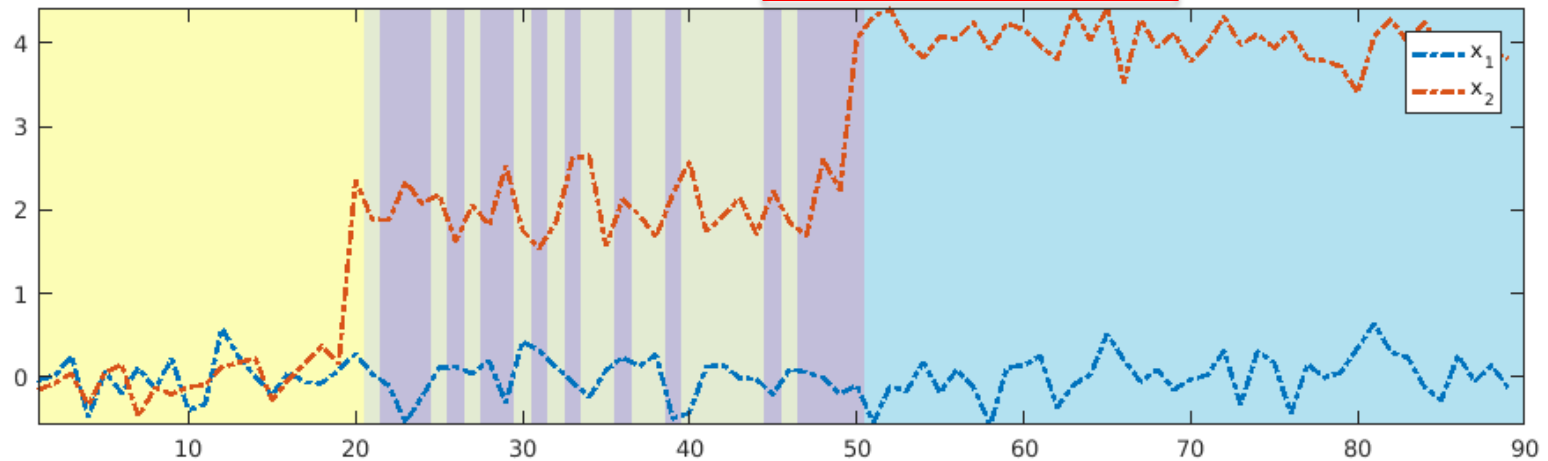
Learning Algorithms and
Systems Laboratory



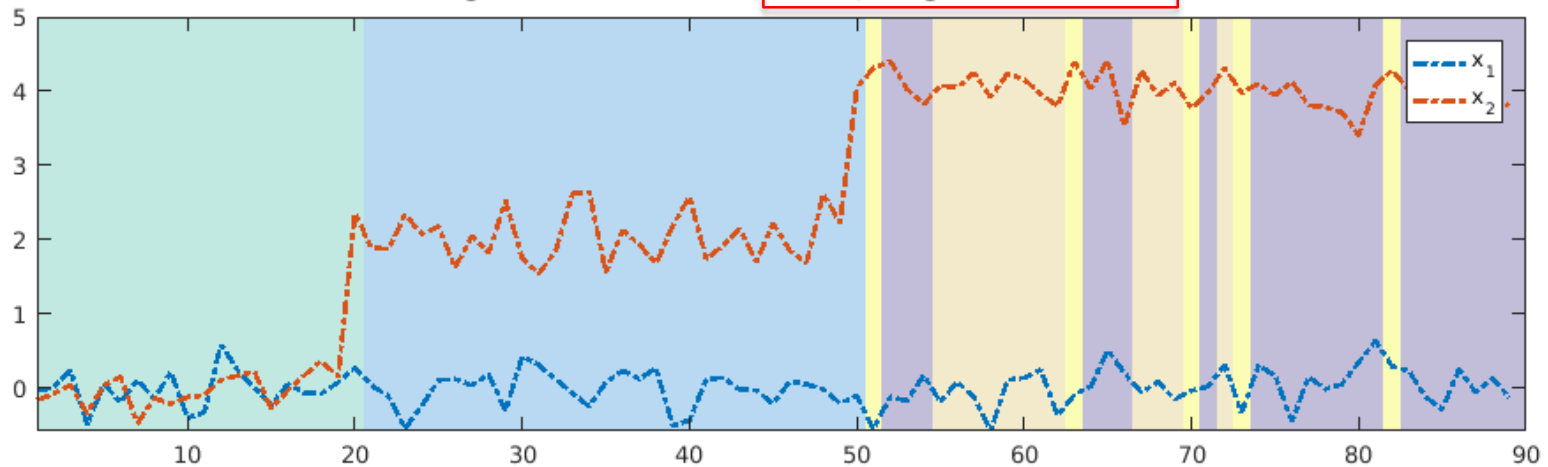
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

HMM for Time series Segmentation

Segmented Data, $K:4$, loglik:-36.7299



Segmented Data, $K:5$, loglik:-49.356



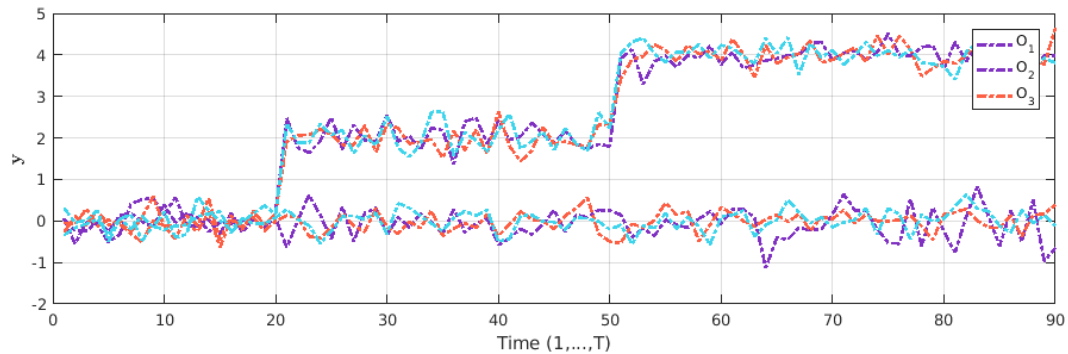
LASA

Learning Algorithms and
Systems Laboratory

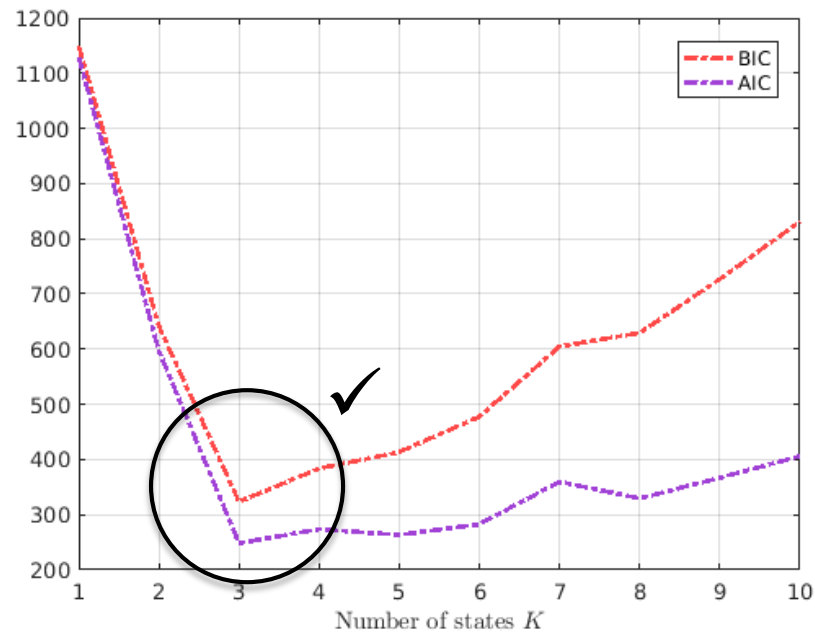


ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

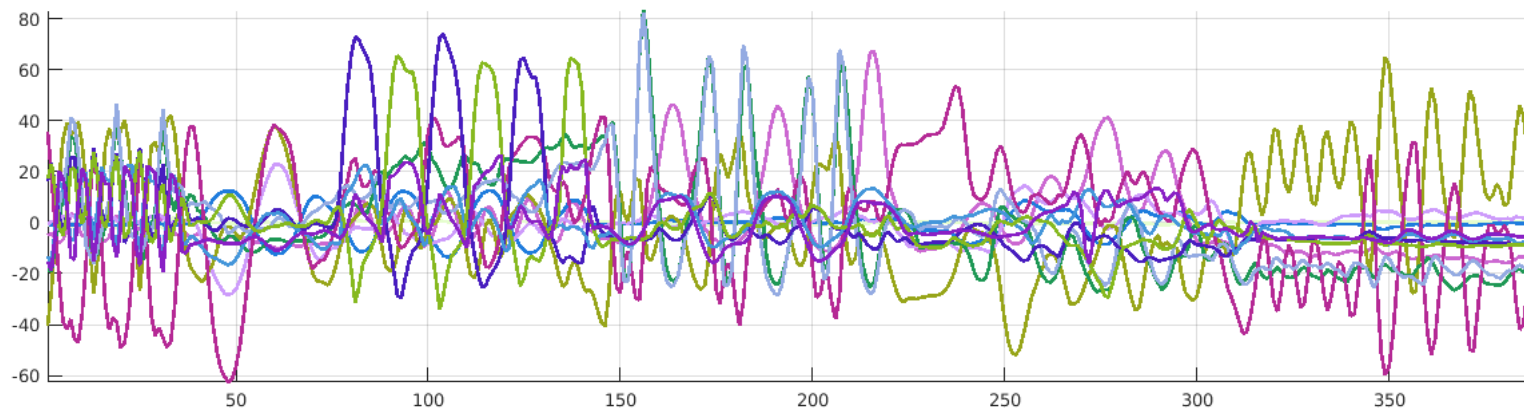
Model Selection for HMMs



HMM Model Selection

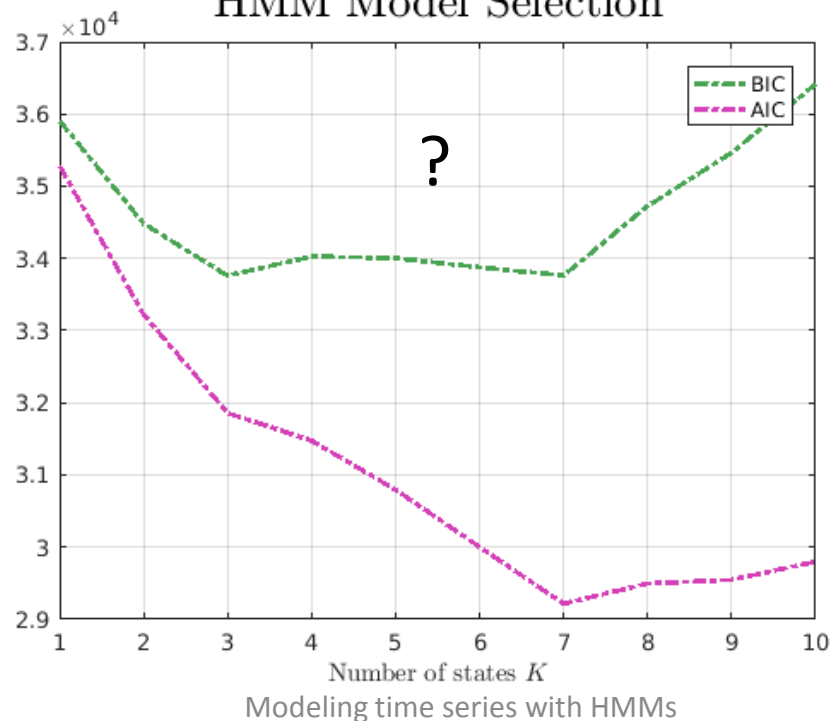


Model Selection for HMMs



$$\mathbf{y}_i \in \mathbb{R}^{12}$$

HMM Model Selection



LASA
Learning Algorithms and
Systems Laboratory



Limitations of classical finite HMMs for Segmentation

Undefined #
hidden states

Cardinality: Choice of hidden states is based on Model Selection **heuristics**, there is **little understanding** of the strengths and weaknesses of such methods in the context of HMMs.

Fixed Transition
Matrix

Topology: We assume that all time series share the **same set** of **emission models** and **switch** among them in exactly the same manner [2].

[1] Emily Fox et al., An HDP-HMM for Systems with State Persistence, ICML, 2008

[2] Emily Fox et al., Sharing Features among Dynamical Systems with Beta Processes, NIPS, 2009

Solution: Bayesian Non-Parametrics

Bayesian Non-Parametrics

- **Bayesian:** Use *Bayesian inference* to estimate the parameters; i.e. **priors** on model **parameters**!

$$\lambda = \{\pi, A, \Theta\}$$

Prior Prior
↓ ↙

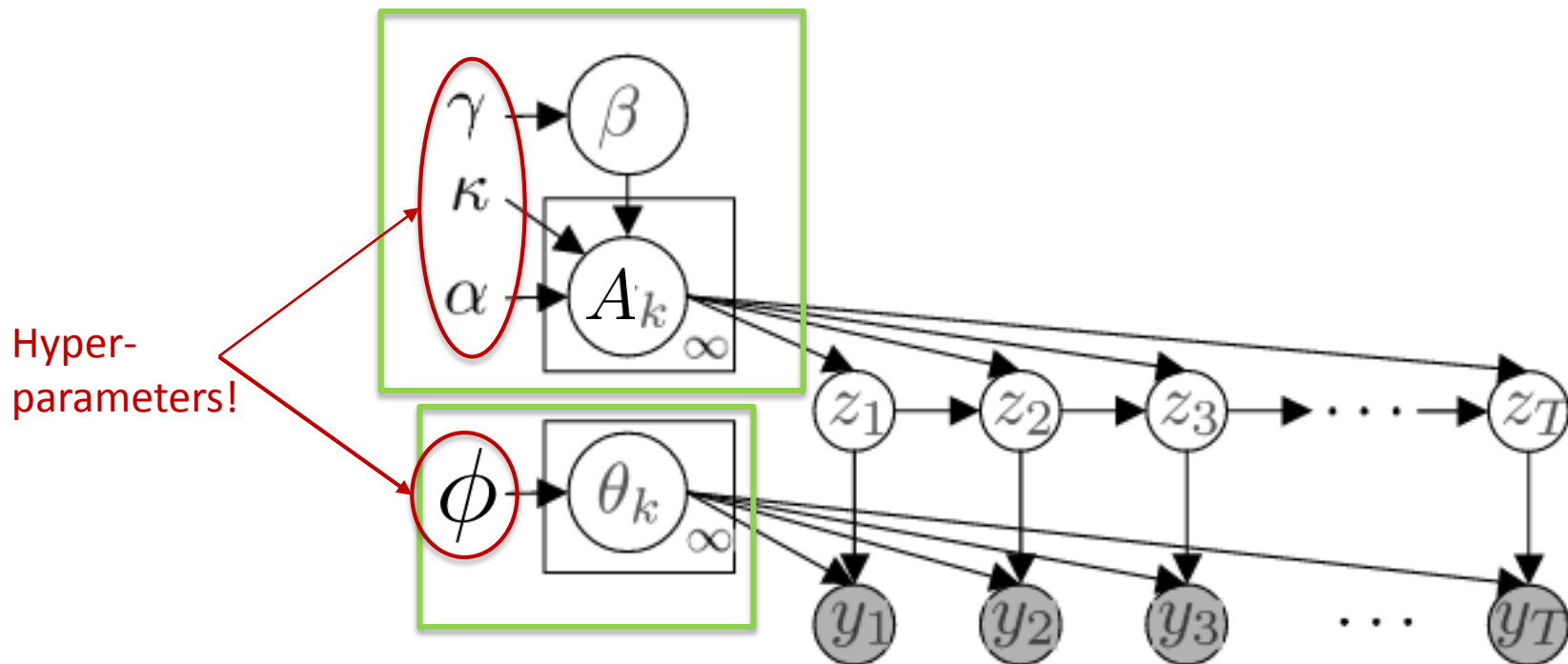
- **Non-parametric:** Does NOT mean methods with “no parameters”, rather models whose complexity (# of states, # Gaussians) is inferred from the data.
 1. Number of parameters grows with sample size.
 2. **Infinite-dimensional** parameter space!



BNP for HMMs: HDP-HMM

☺ Cardinality

- **Hierarchical Dirichlet Process (HDP)** prior on the transition Matrix!



- **Normal Inverse Wishart (NIW)** prior on emission parameters!

Emily Fox et al., An HDP-HMM for Systems with State Persistence, ICML, 2008

BNP for HMMs: HDP-HMM

☺ Cardinality

- The **Dirichlet Process (DP)** is a prior distribution over distributions.

$$M(\cdot) = \sum_{k=1}^{\infty} C_k \delta_{\Theta_k}(\cdot) \quad \text{where} \quad \sum_{k=1}^{\infty} C_k = 1$$

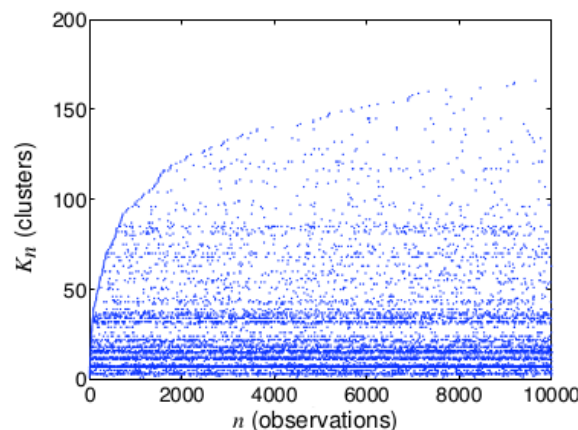
- Used for clustering with infinite mixture models; i.e. instead of setting K in a GMM, the K is learned from data.

$K_n = \#$ clusters in sample of size n

This only gives us an estimate of the K clusters!

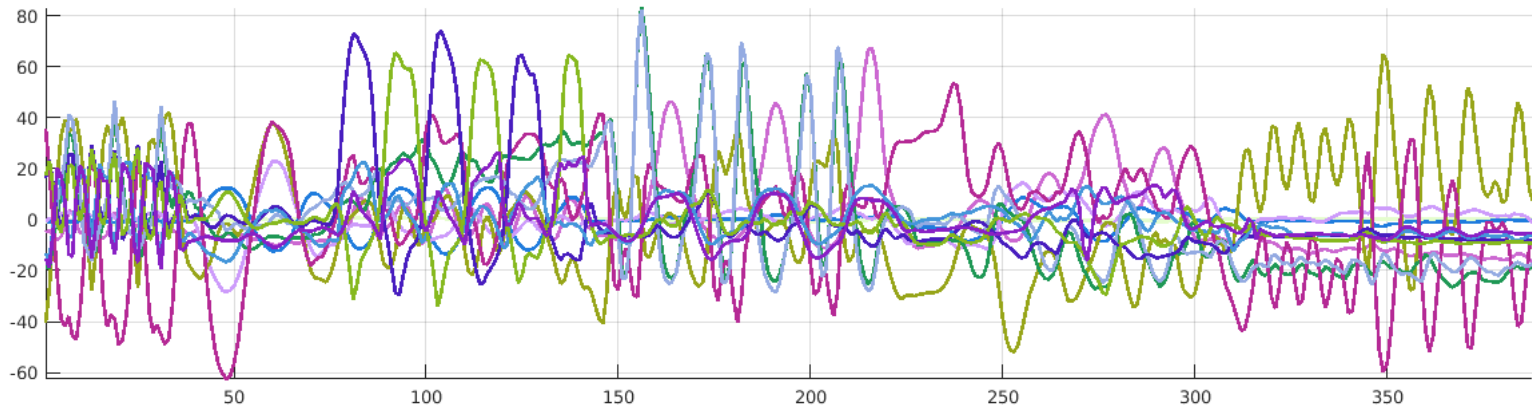
Cannot be use directly on transition matrix:

$$A \in \mathbb{R}^{K \times K}$$

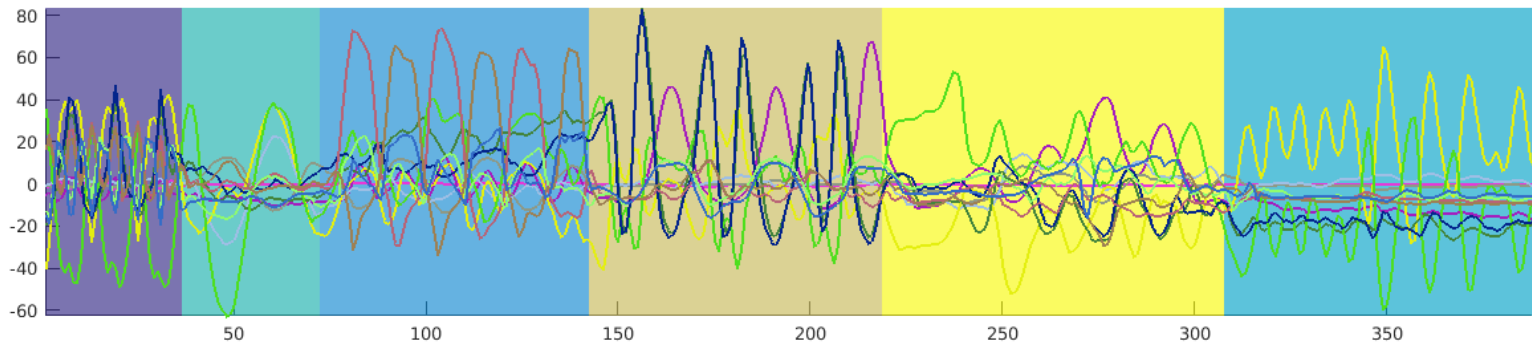


- The **Hierarchical Dirichlet Process (HDP)** is a hierarchy of DPs!

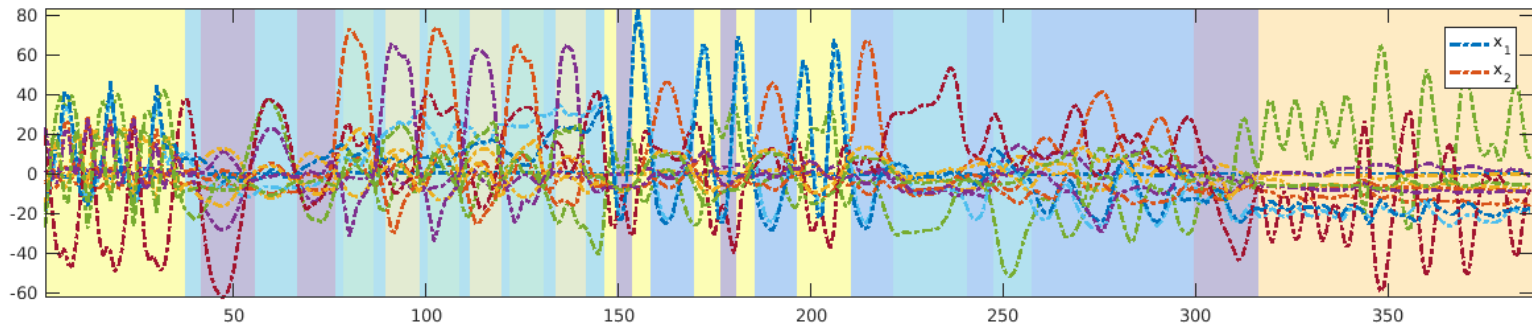
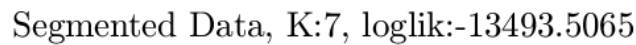
Segmentation with HDP-HMM



$$\mathbf{y}_i \in \mathbb{R}^{12}$$



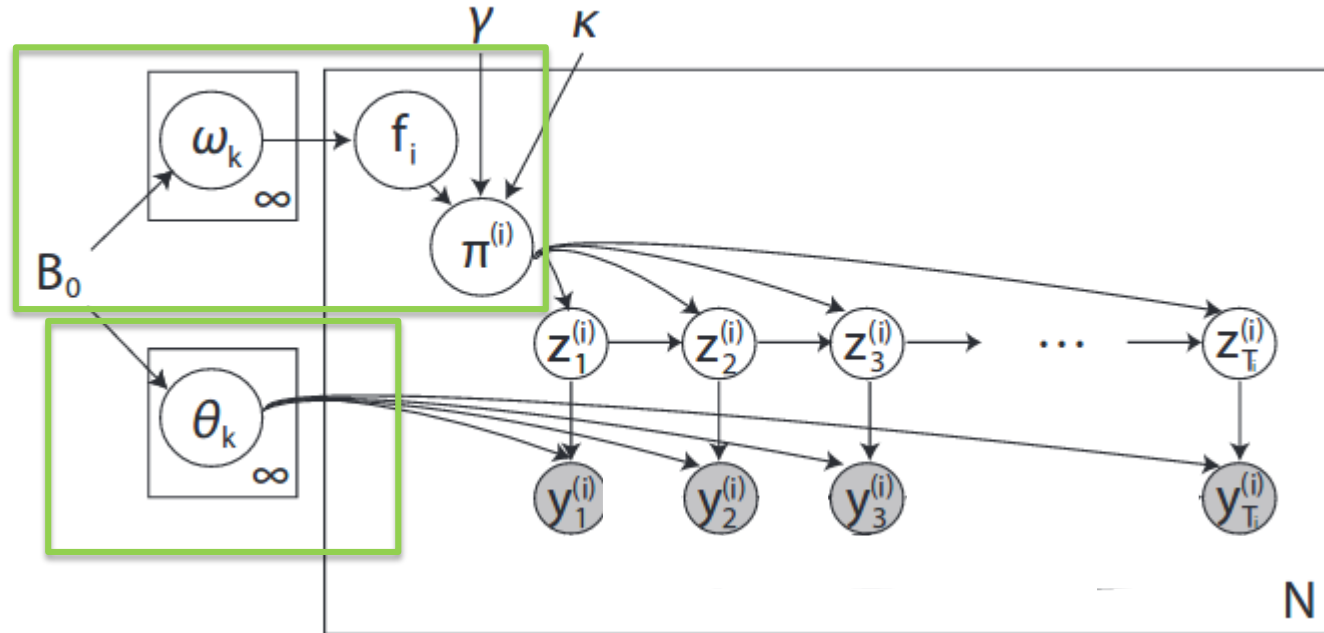
Emily Fox et al., An HDP-HMM for Systems with State Persistence, ICML, 2008

$$\mathbf{y}_i \in \mathbb{R}^{12}$$


BNP for HMMs: BP-HMM

☺ Cardinality
☺ Topology

- The **Beta Process (BP)** prior on the transition Matrix!



- **Normal Inverse Wishart (NIW)** prior on emission parameters!

Emily Fox et al., Sharing Features among Dynamical Systems with Beta Processes, NIPS, 2009

BNP for HMMs: BP-HMM

☺ Cardinality
☺ Topology

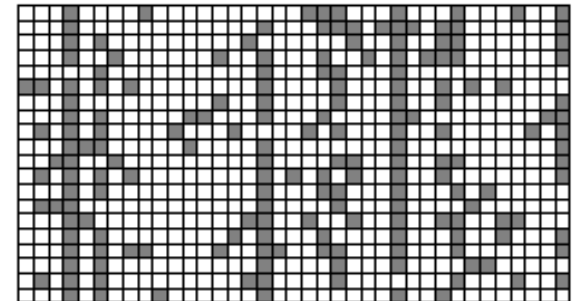
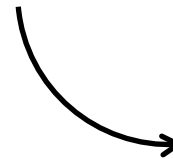
- The **Beta Process (BP)** prior on the transition Matrix!

Beta Process (BP)

Distribution on objects of the form

$$\theta = \sum_{k=1}^{\infty} w_k \delta_{\phi_k} \quad \text{with } w_k \in [0, 1] .$$

Time-Series



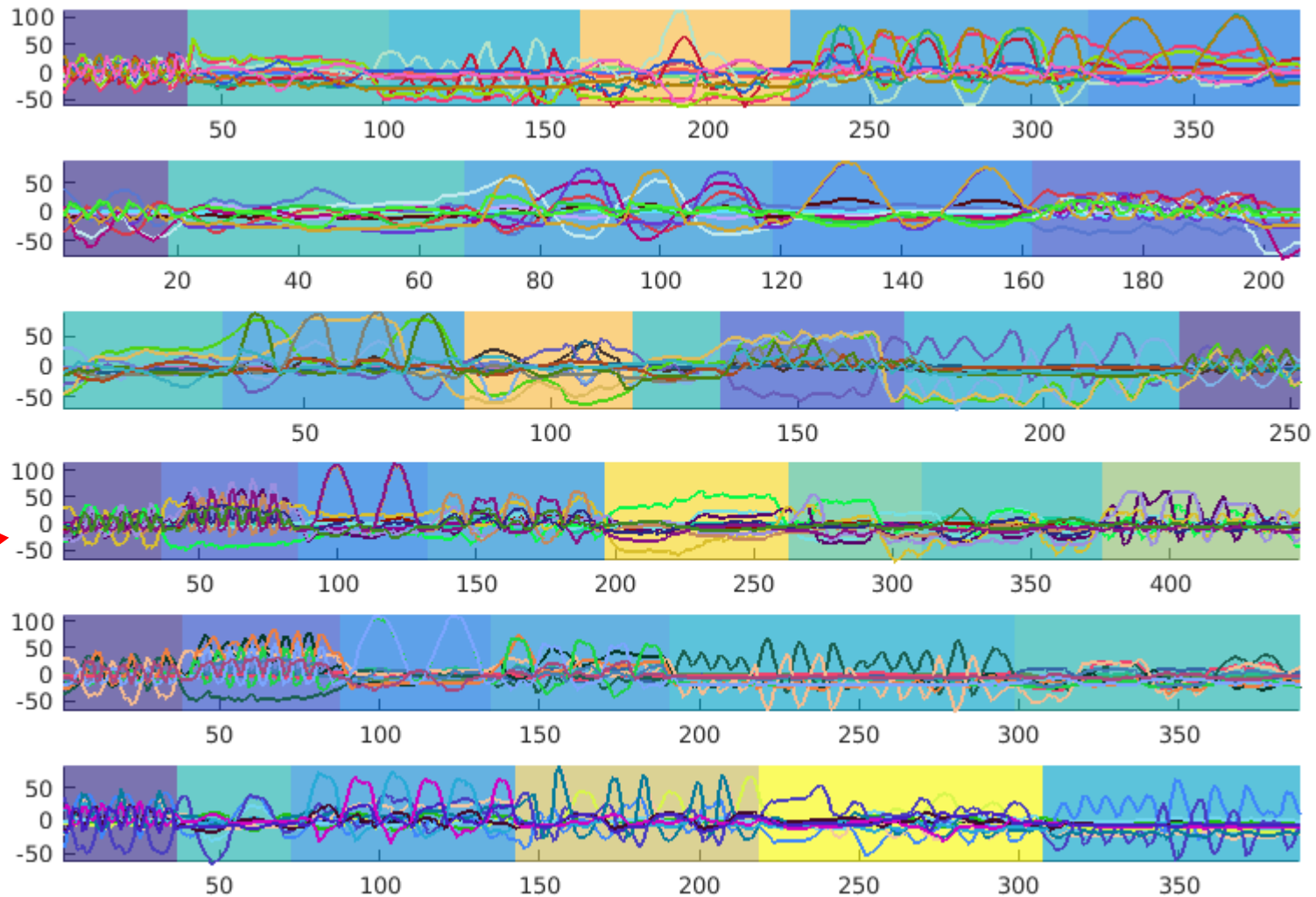
Features (i.e. shared HMM States)

Segmentation in Human Motion Data

12 Variables

- Torso position
- Waist Angles (2)
- Neck Angle
- Shoulder Angles
- ..

$$\mathbf{y}_i \in \mathbb{R}^{12}$$



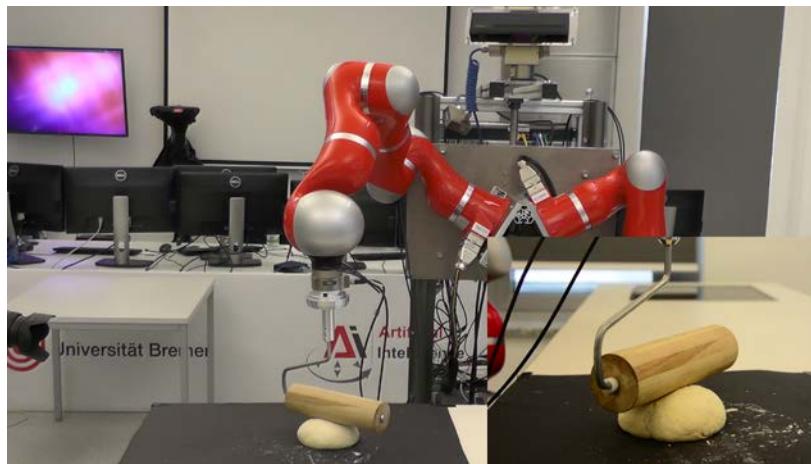
Emily Fox et al., Sharing Features among Dynamical Systems with Beta Processes, NIPS, 2009



Applications in Robotics



Learning Complex Sequential Tasks from Demonstration



Learning Complex Sequential Tasks from Demonstrations: Pizza Dough Rolling

Nadia Figueroa, Lucia Pais and Aude Billard



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE



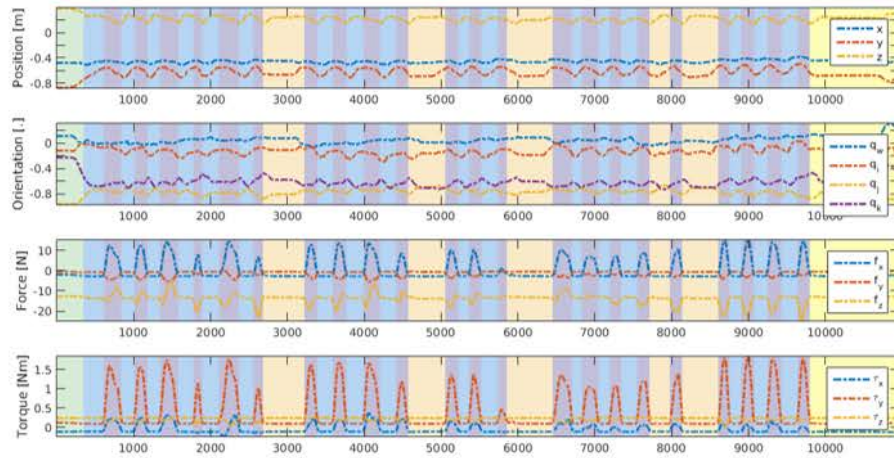
LASA

Learning Algorithms and
Systems Laboratory

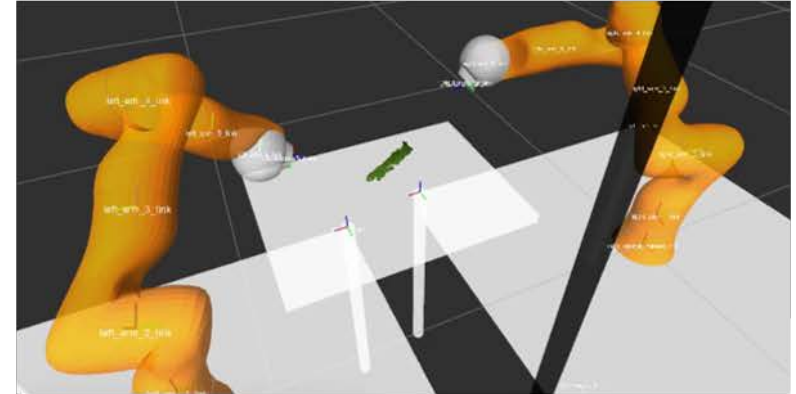
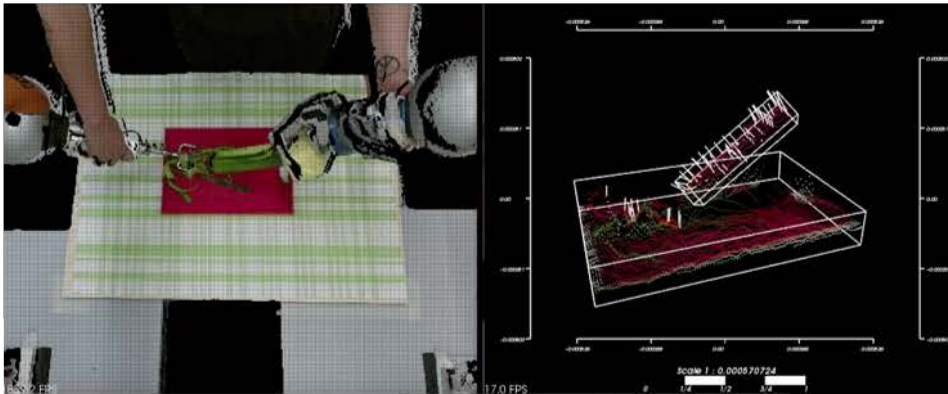
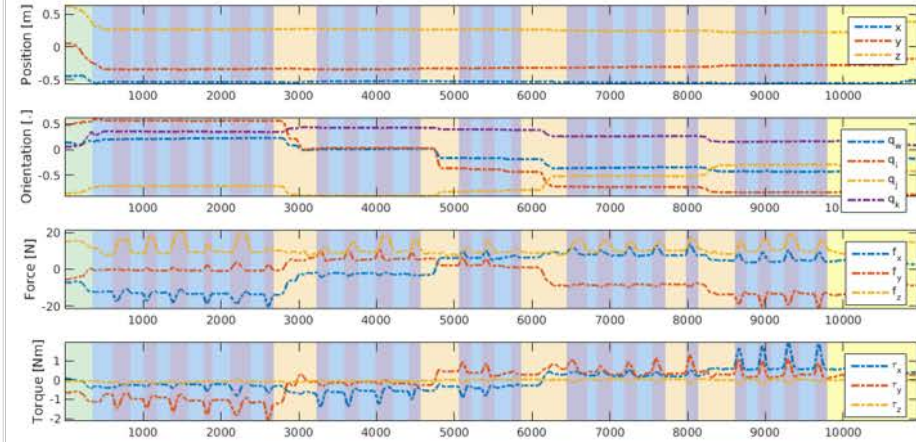


ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

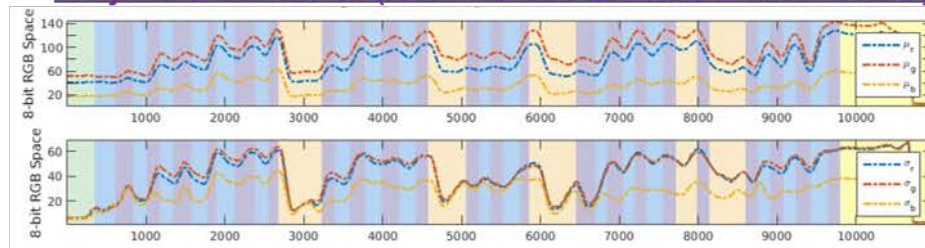
Right Arm EE data (pos,orientation,wrench)

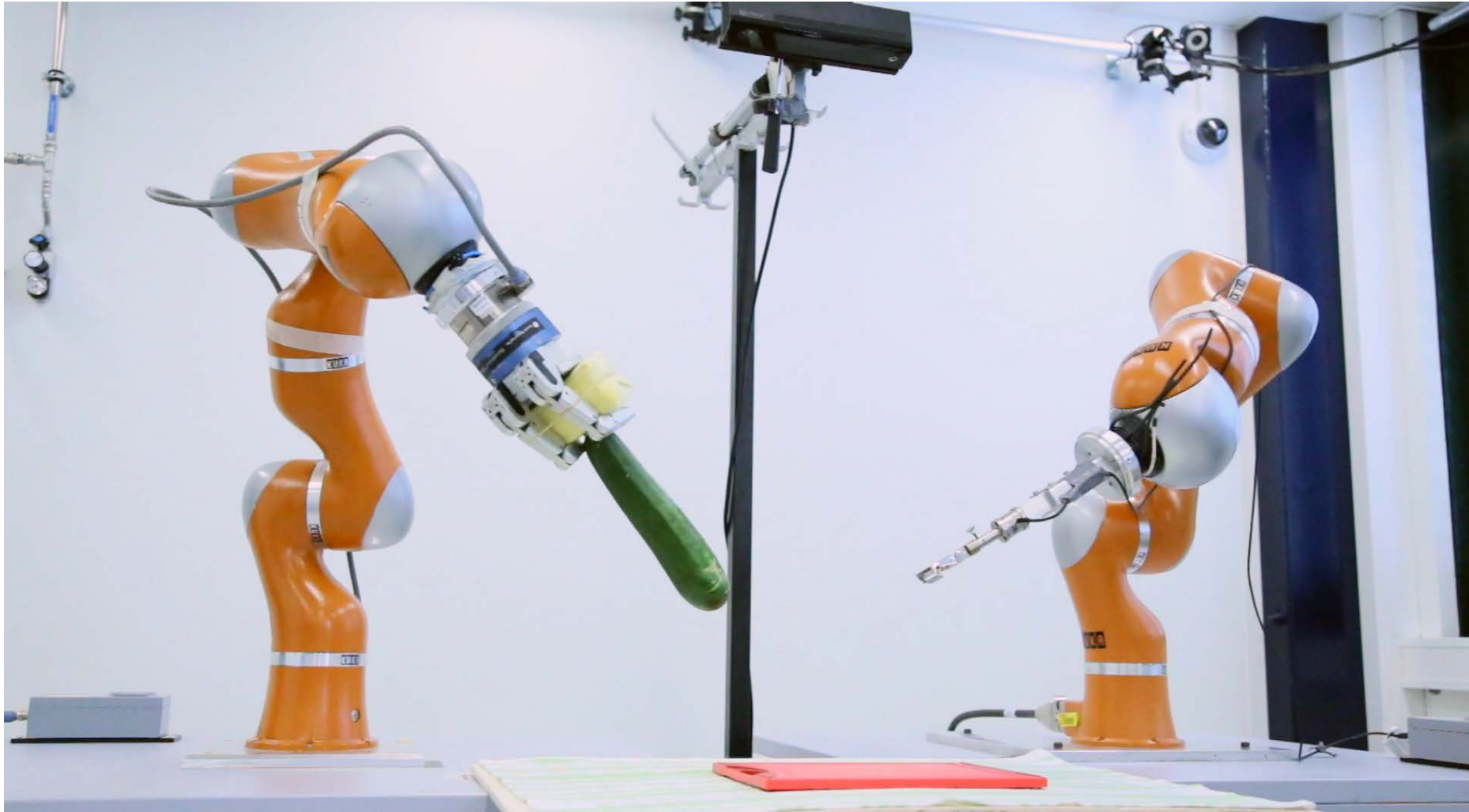


Left Arm EE data (pos,orientation,wrench)



Object Features (mean, std of RGB Channels)





LASA
Learning Algorithms and
Systems Laboratory



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE