

EE613 - Machine Learning for Engineers

HIDDEN MARKOV MODELS

Sylvain Calinon

Robot Learning and Interaction Group

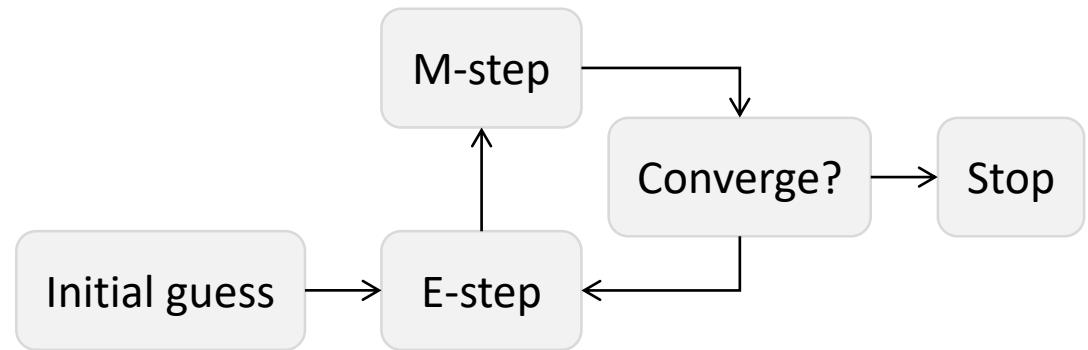
Idiap Research Institute

Nov 11, 2021

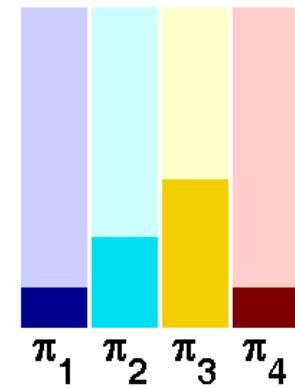
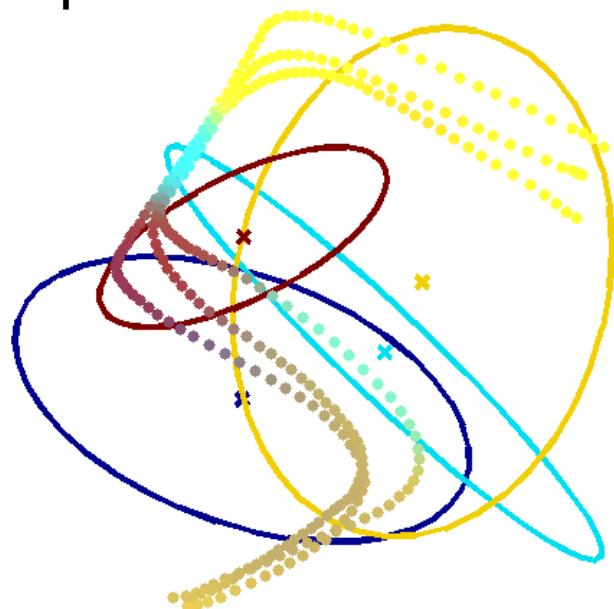
Outline

- Markov models
- Hidden Markov model (HMM)
- Forward-backward algorithm
- Viterbi decoding (dynamic programming)
- Hidden semi-Markov model (HSMM)
- HMM with dynamic features (Trajectory-HMM)

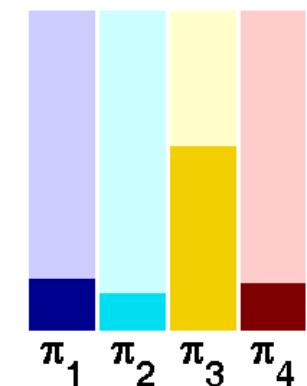
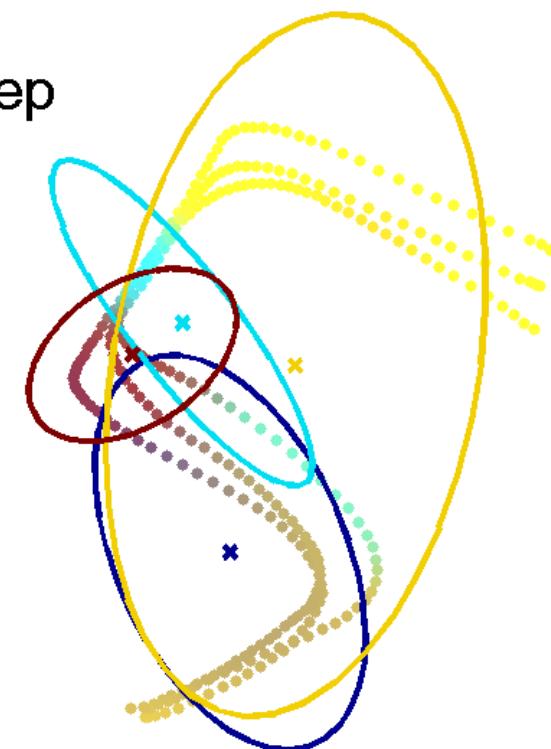
Recap: EM for GMM



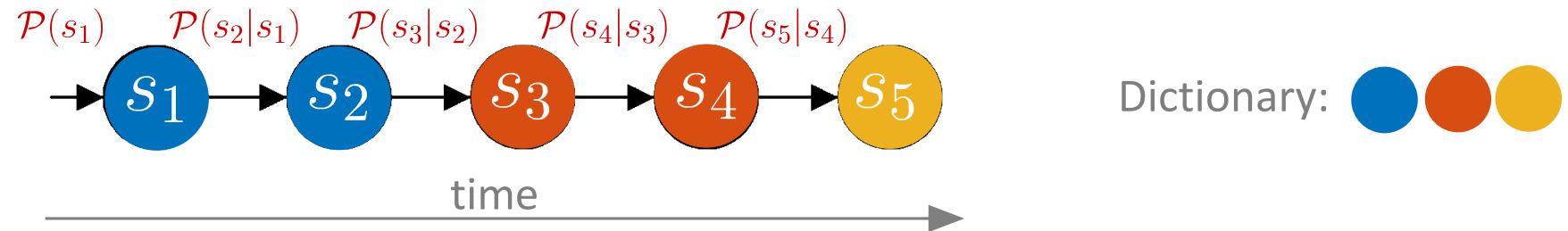
E-step



M-step



Markov models



With a **first order Markov model**, the joint distribution of a sequence of states is assumed to be of the form

$$\mathcal{P}(s_1, s_2, \dots, s_T) = \mathcal{P}(s_1) \prod_{t=2}^T \mathcal{P}(s_t | s_{t-1})$$

and we thus have

$$\mathcal{P}(s_t | s_1, s_2, \dots, s_{t-1}) = \mathcal{P}(s_t | s_{t-1})$$

In most applications, the conditional distributions $\mathcal{P}(s_t | s_{t-1})$ will be assumed to be **stationary (homogeneous Markov chain)**.

Markov models – Parameters

K possible states

The **initial state distribution** is defined by

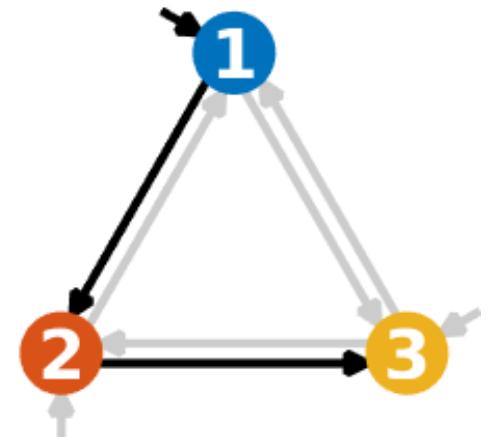
$$\Pi_i = \mathcal{P}(s_1=i) \quad \text{with} \quad \sum_{i=1}^K \Pi_i = 1$$

A **transition matrix A** is defined, with elements

$$a_{i,j} = \mathcal{P}(s_{t+1}=j \mid s_t=i)$$

defining the probability of getting from state i to state j in one step.

Constraint: each row of the matrix sums to one, $\sum_{j=1}^K a_{i,j} = 1$.



	1	2	3
1	$a_{1,1}$	$a_{1,2}$	$a_{1,3}$
2	$a_{2,1}$	$a_{2,2}$	$a_{2,3}$
3	$a_{3,1}$	$a_{3,2}$	$a_{3,3}$

Example: language modeling

We define the state space to be all the words in some language.

The marginal probabilities $\mathcal{P}(s_t = k)$ are called **unigram** statistics.

For a first-order Markov model, $\mathcal{P}(s_t = k \mid s_{t-1} = j)$ is called a **bigram** model.

For a second-order Markov model, $\mathcal{P}(s_t = k \mid s_{t-1} = j, s_{t-2} = i)$ is called a **trigram** model, etc.

In the general case, these are called **n -gram** models.

Example: language modeling

Sentence completion

The model can predict the next word given the previous words in a sentence. This can be used to reduce the amount of typing required (e.g., mobile devices).

Data compression

The model can be used to define an encoding scheme, by assigning codewords to more probable strings. The more accurate the predictive model, the fewer the number of bits is required to store the data.

Text classification

The model can be used as a class-conditional density and/or generative classifier.

Automatic writing

The model can be used to sample from $\mathcal{P}(s_1, s_2, \dots, s_t)$ to generate artificial text.

Example: language modeling

SAYS IT'S NOT IN THE CARDS LEGENDARY RECONNAISSANCE BY
ROLLIE DEMOCRACIES UNSUSTAINABLE COULD STRIKE
REDLINING VISITS TO PROFIT BOOKING WAIT HERE AT
MADISON SQUARE GARDEN COUNTY COURTHOUSE WHERE HE
HAD BEEN DONE IN THREE ALREADY IN ANY WAY IN WHICH A
TEACHER ...

Example of text generated from a 4-gram model, trained on a corpus of 400 million words.

The first 4 words are specified by hand, the model generates the 5th word, and then the results are fed back into the model.

Source: <http://www.fit.vutbr.cz/~imikolov/rnnlm/gen-4gram.txt>

MLE of transition matrix in Markov models

A Markov model is described by $\Theta^{\text{MM}} = \left\{ \{a_{i,j}\}_{j=1}^K, \Pi_i \right\}_{i=1}^K$, where the transition probabilities $a_{i,j}$ are stored in a matrix \mathbf{A} .

The maximum likelihood estimate (MLE) of the parameters can be computed with the normalized counts

$$\hat{\Pi}_i = \frac{N_i}{\sum_{k=1}^K N_k}, \quad \hat{a}_{i,j} = \frac{N_{i,j}}{\sum_{k=1}^K N_{i,k}}$$

These results can be extended to higher order Markov models, but since an n-gram models has $O(K^n)$ parameters, special care needs to be taken with overfitting.

For example, with a bi-gram model and 50,000 words in the dictionary, there are 2.5 billion parameters to estimate, and it is unlikely that all possible transitions will be observed in the training data.

Hidden Markov model (HMM)

Python notebook:
`demo_HMM.ipynb`

Matlab code:
`demo_HMM01.m`

Hidden Markov model (HMM)

In a Markov chain, the state is directly visible to the observer → the transition probabilities are the only parameters.

In an HMM, the state is not directly visible, but an output dependent on the state is visible.

Initial state: All we know is that 60% of the days are rainy on average.

You can think of an HMM either as:

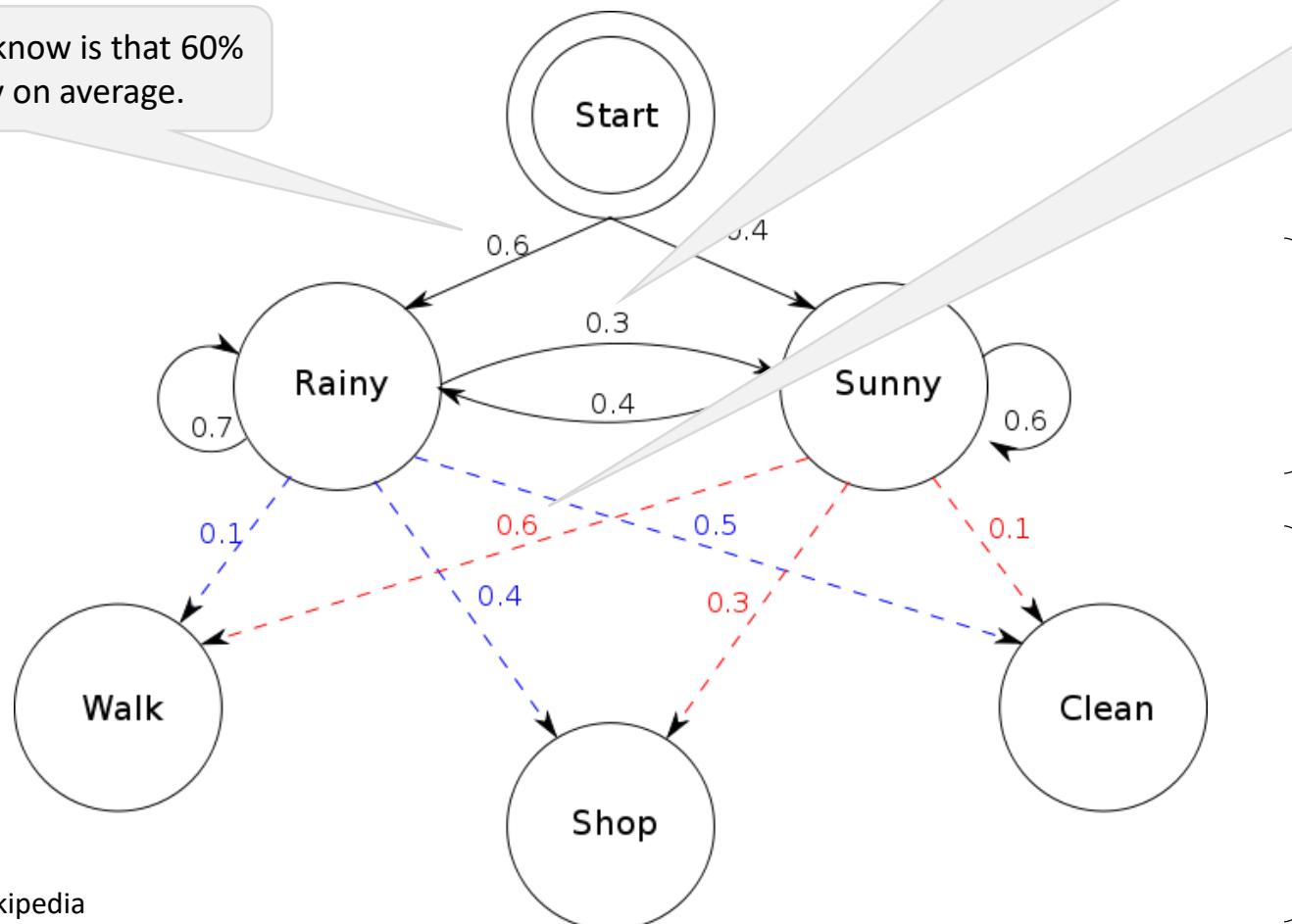
- a Markov chain with stochastic measurements
- a GMM with latent variables changing over time

The **transition probability** represents the change of the weather in the underlying Markov chain.

Here, there is a 30% chance that tomorrow will be sunny if today is rainy.

The **emission probability** represents how likely Bob performs a certain activity on each day.
If it is sunny, there is a 60% chance that he is outside for a walk.

If it is rainy, there is a 50% chance that he cleans his apartment, etc.



Hidden states

Observed output
(emission probability)

Inference problems associated with HMMs

Probability of an observed sequence

$$\mathcal{P}(\boldsymbol{\xi}_{1:T}) = \mathcal{P}(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_T)$$

Probability of the latent variables

Filtering

$$\mathcal{P}(s_t | \boldsymbol{\xi}_{1:t}) = \mathcal{P}(s_t | \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_t)$$

Prediction

$$\mathcal{P}(s_{t+1} | \boldsymbol{\xi}_{1:t}) = \mathcal{P}(s_{t+1} | \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_t)$$

Smoothing → *Forward-backward algorithm*

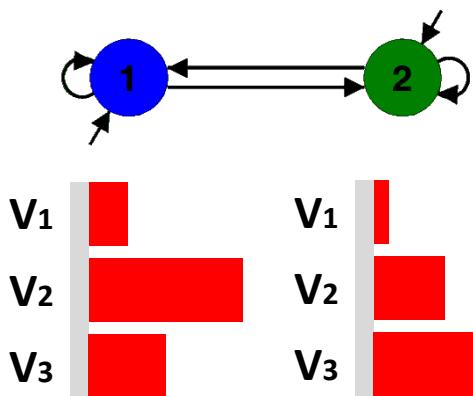
$$\mathcal{P}(s_t | \boldsymbol{\xi}_{1:T}) = \mathcal{P}(s_t | \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_T)$$

MAP estimation → *Viterbi decoding*

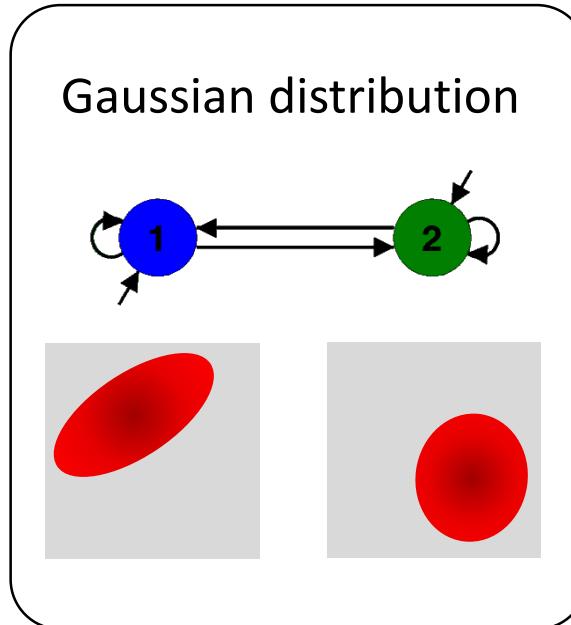
$$\mathcal{P}(s_{1:T} | \boldsymbol{\xi}_{1:T}) = \mathcal{P}(s_1, s_2, \dots, s_T | \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_T)$$

Emission/output distributions in HMM

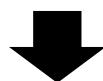
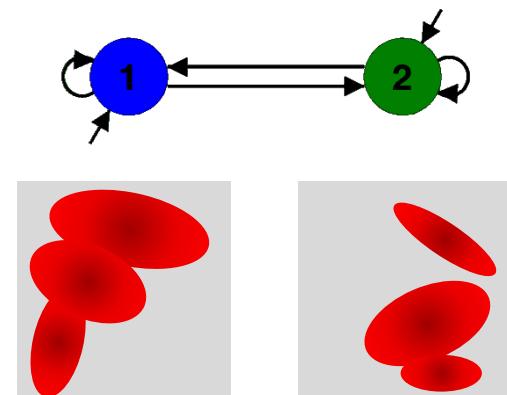
Discrete tables



Gaussian distribution

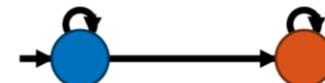
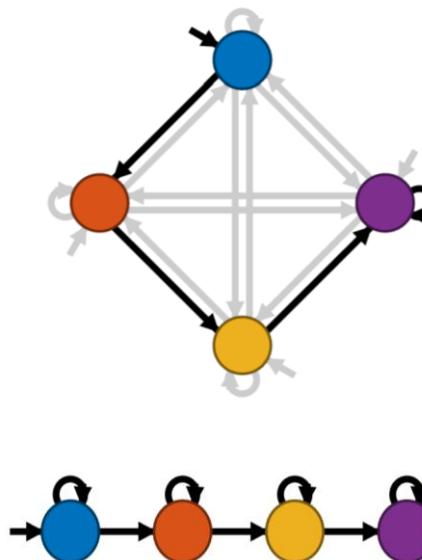
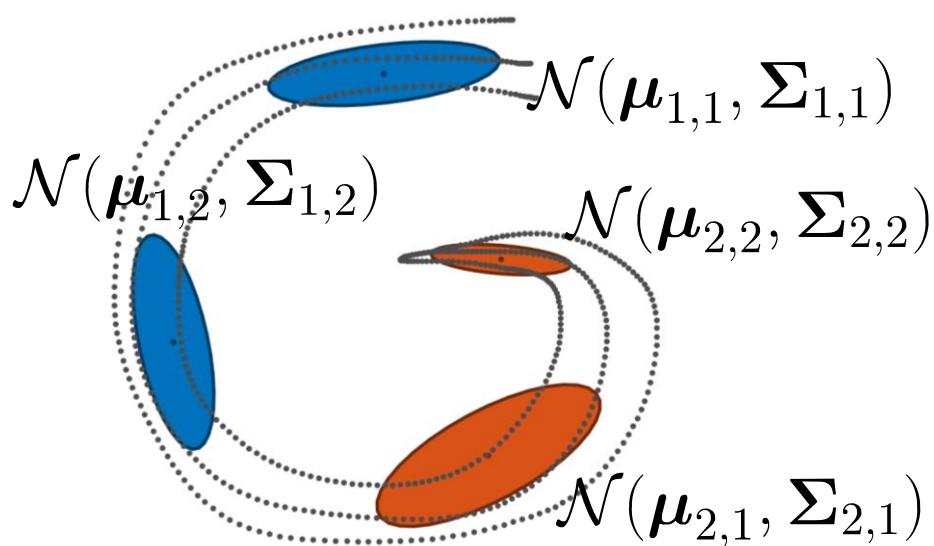
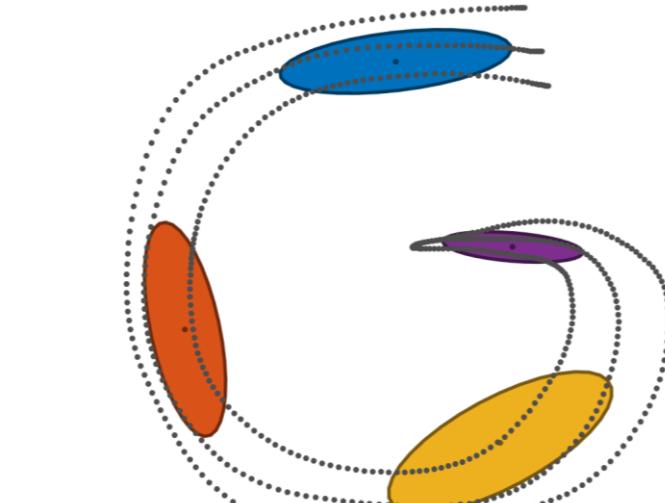


Mixture of Gaussians



GMM with latent variable z_t depending
on the conditional distribution $\mathcal{P}(z_t|z_{t-1})$

Transition matrix structures in HMM



$$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{bmatrix}$$

$$\begin{bmatrix} a_{1,1} & a_{1,2} & 0 & 0 \\ 0 & a_{2,2} & a_{2,3} & 0 \\ 0 & 0 & a_{3,3} & a_{3,4} \\ 0 & 0 & 0 & a_{4,4} \end{bmatrix}$$

$$\begin{bmatrix} a_{1,1} & a_{1,2} \\ 0 & a_{2,2} \end{bmatrix}$$

HMM - Examples of application

HMM is used in many fields as a tool for **time series or sequences analysis**, and in fields where the goal is to recover a data sequence that is not immediately observable:

- Speech recognition
 - Speech synthesis
 - Part-of-speech tagging
 - Natural language modeling
 - Machine translation
 - Gene prediction
 - Molecule kinetic analysis
 - DNA motif discovery
 - Alignment of bio-sequences (e.g., proteins)
 - Metamorphic virus detection
 - Document separation in scanning solutions
 - Cryptoanalysis
 - Activity recognition
 - Protein folding
 - Human motion science
 - Online handwriting recognition
 - Robotics
- and many,
many others...*

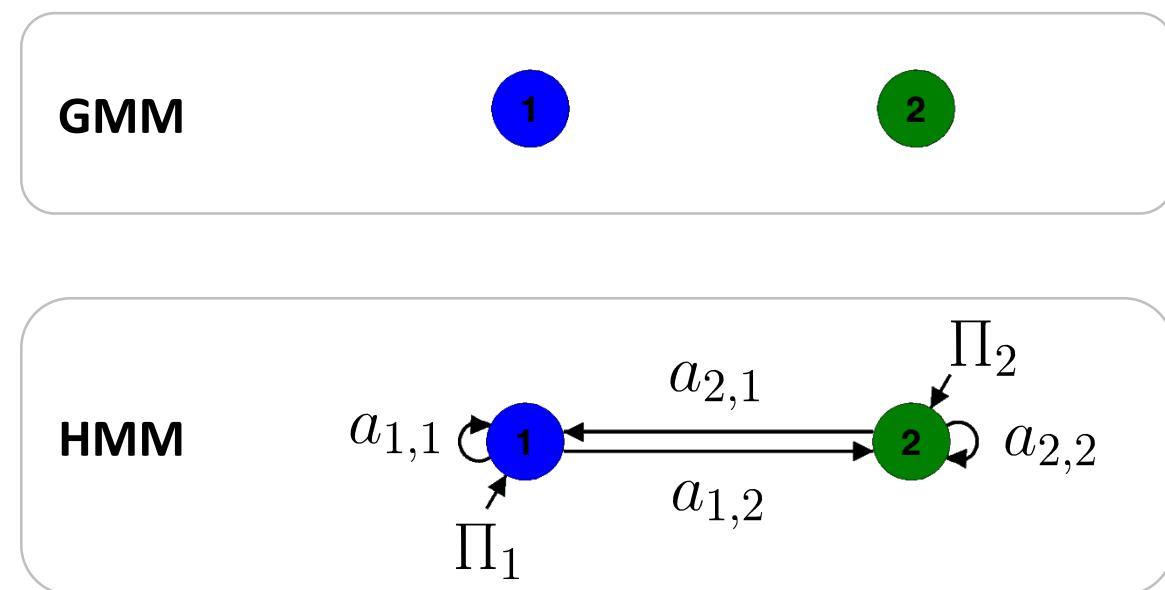
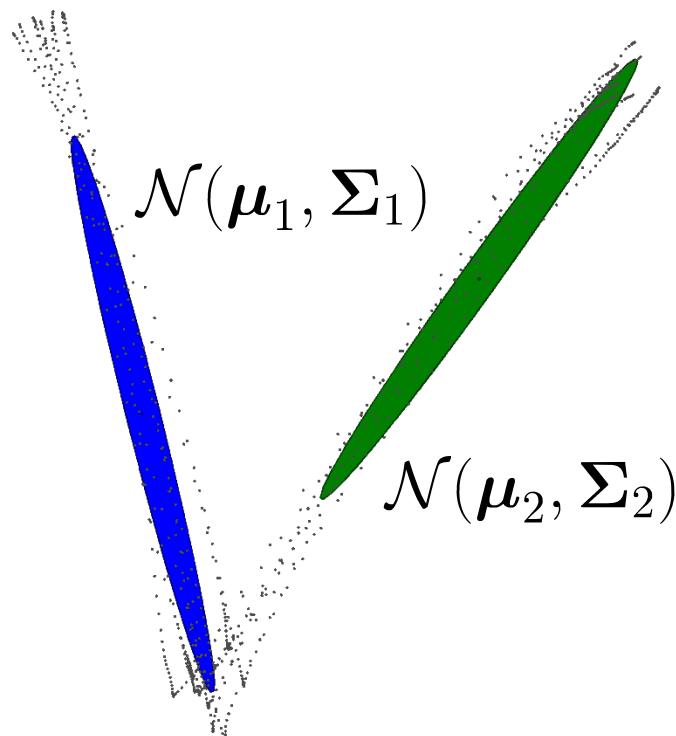
HMM parameters

$$\Theta^{\text{GMM}} = \{\pi_i, \mu_i, \Sigma_i\}_{i=1}^K$$

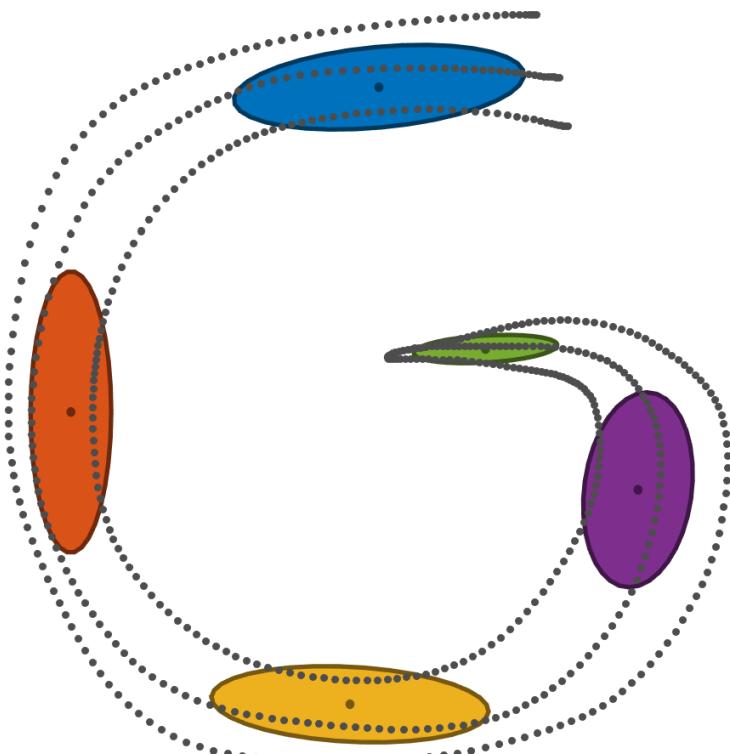
$$\Theta^{\text{HMM}} = \{\{a_{i,j}\}_{j=1}^K, \Pi_i, \mu_i, \Sigma_i\}_{i=1}^K$$

From now on, we will consider
a single Gaussian as state output

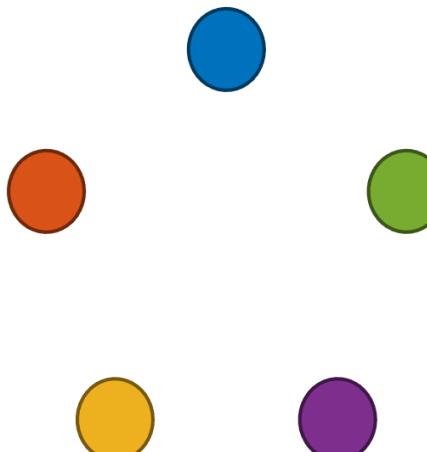
$$\pi_i = 1$$

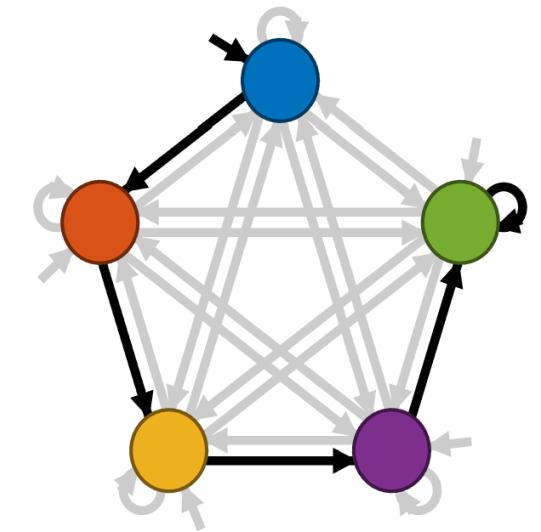
GMM vs HMM



GMM



HMM



Useful intermediary variables in HMM

Forward variable

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$$

Backward variable

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} \mid s_t=i)$$

Smoothed node marginals

$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i \mid \xi_{1:T})$$

Smoothed edge marginals

$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j \mid \xi_{1:T})$$

Forward algorithm

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \boldsymbol{\xi}_{1:t})$$

The probability to be in state i at time step t given the partial observation $\boldsymbol{\xi}_{1:t} = \{\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_t\}$ can be computed with the **forward variable**

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_t) = \mathcal{P}(s_t=i, \boldsymbol{\xi}_{1:t})$$

which can be used to compute

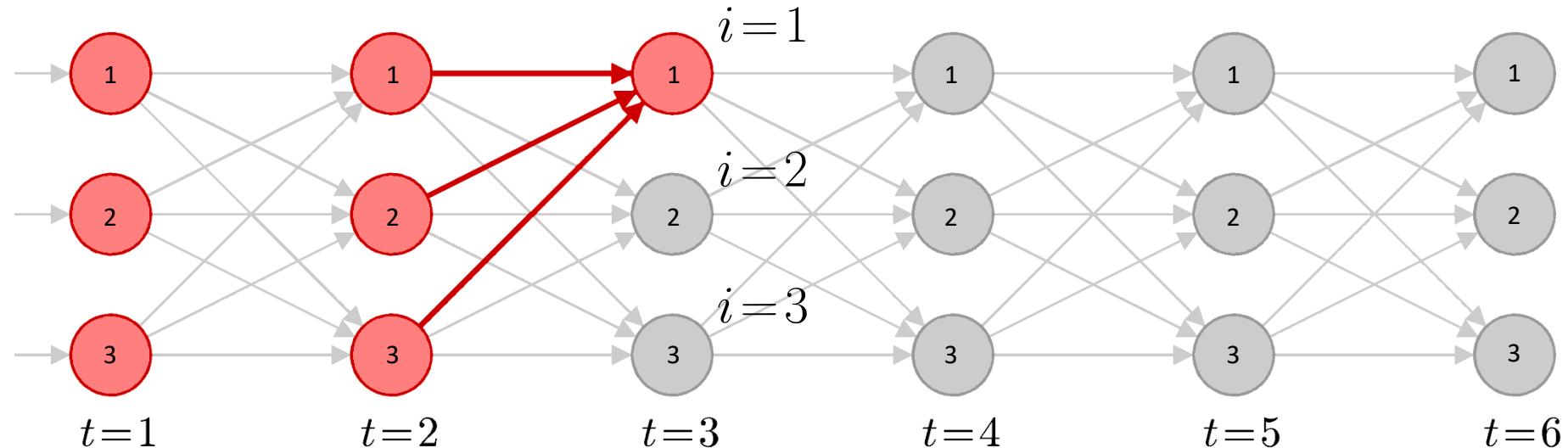
$$\mathcal{P}(s_t=i \mid \boldsymbol{\xi}_{1:t}) = \frac{\mathcal{P}(s_t=i, \boldsymbol{\xi}_{1:t})}{\mathcal{P}(\boldsymbol{\xi}_{1:t})} = \frac{\alpha_{t,i}^{\text{HMM}}}{\sum_{k=1}^K \alpha_{t,k}^{\text{HMM}}}$$

The direct computation would require marginalizing over all possible state sequences $\{s_1, s_2, \dots, s_{t-1}\}$, which would grow exponentially with t .

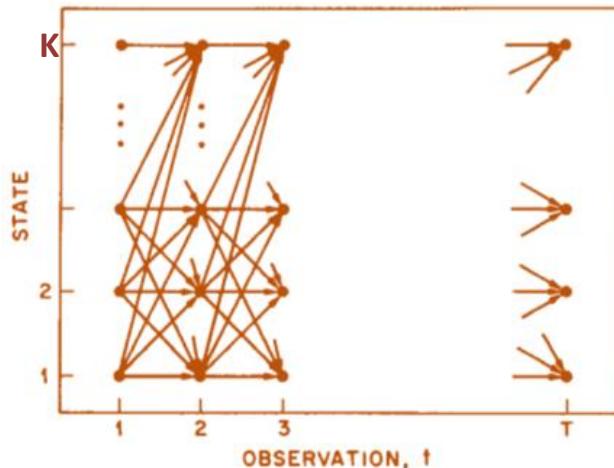
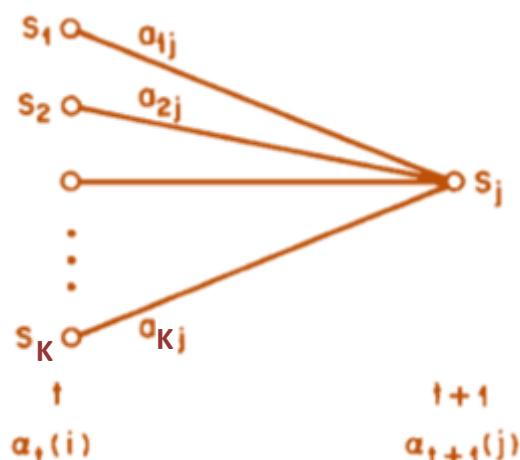
The forward algorithm takes advantage of the conditional independence rules of the HMM to perform the calculation recursively.

Forward algorithm

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$$



$$\alpha_{t,i}^{\text{HMM}} = \left(\sum_{j=1}^K \alpha_{t-1,j}^{\text{HMM}} a_{j,i} \right) \mathcal{N}(\xi_t | \mu_i, \Sigma_i) \quad \text{with} \quad \alpha_{1,i}^{\text{HMM}} = \prod_i \mathcal{N}(\xi_1 | \mu_i, \Sigma_i)$$



It can be used to evaluate trajectories by computing the likelihood

$$\mathcal{P}(\xi | \Theta^{\text{HMM}}) = \sum_{i=1}^K \alpha_{T,i}^{\text{HMM}}$$

Useful intermediary variables in HMM

Forward variable

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$$

Backward variable

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} \mid s_t=i)$$

Smoothed node marginals

$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i \mid \xi_{1:T})$$

Smoothed edge marginals

$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j \mid \xi_{1:T})$$

Backward algorithm

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\boldsymbol{\xi}_{t+1:T} \mid s_t = i)$$

Similarly, we can define a **backward variable** starting from the boundary condition

$$\beta_{T,i}^{\text{HMM}} = 1$$

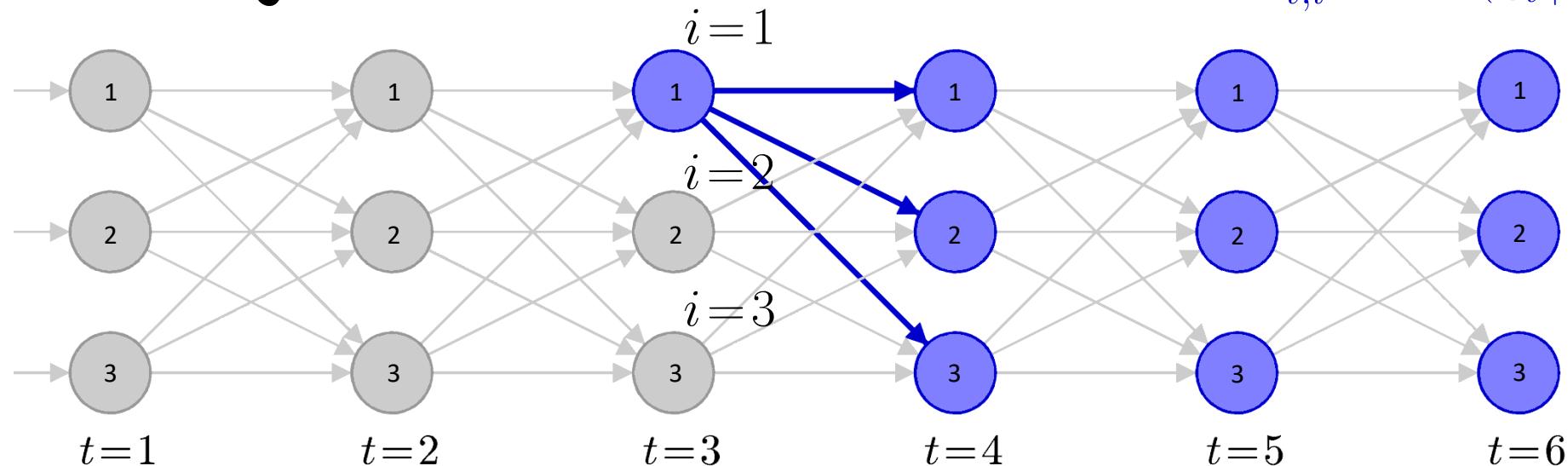
and computed as

$$\beta_{t,i}^{\text{HMM}} = \sum_{j=1}^K a_{i,j} \mathcal{N}(\boldsymbol{\xi}_{t+1} \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \beta_{t+1,j}^{\text{HMM}}$$

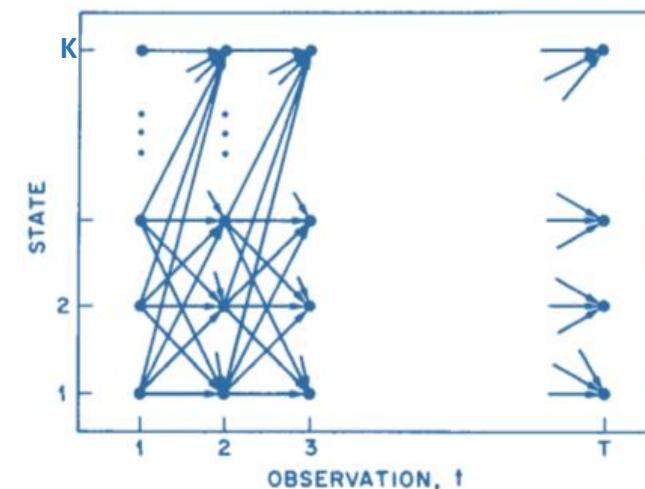
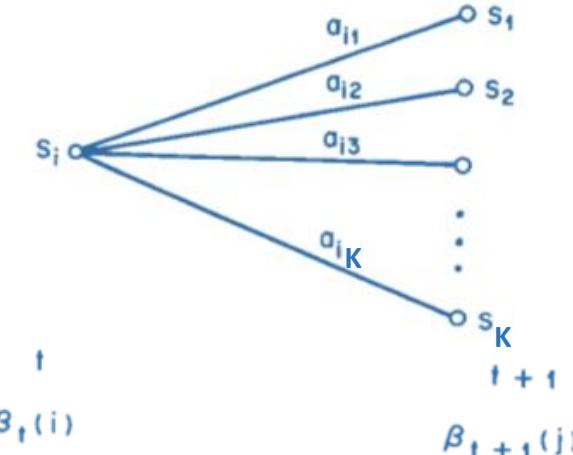
corresponding to the probability of the partial observation $\{\boldsymbol{\xi}_{t+1}, \dots, \boldsymbol{\xi}_{T-1}, \boldsymbol{\xi}_T\}$, knowing that we are in state i at time step t .

Backward algorithm

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} \mid s_t = i)$$



$$\beta_{t,i}^{\text{HMM}} = \sum_{j=1}^K a_{i,j} \mathcal{N}(\xi_{t+1} \mid \mu_j, \Sigma_j) \beta_{t+1,j}^{\text{HMM}} \quad \text{with} \quad \beta_{T,i}^{\text{HMM}} = 1$$



boundary condition,
corresponding to
 $\mathcal{P}(\emptyset \mid s_T = i) = 1, \forall i$

Useful intermediary variables in HMM

Forward variable

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{F}$$

These variables are sometimes called ***smoothed variables*** as they combine forward and backward probabilities in the computation.

You can think of their roles as passing "messages" from left to right, and from right to left, and then combining the information at each node.

Backward variable

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\boldsymbol{\xi}_{t+1:T} | s_t)$$

Smoothed node marginals

$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i | \boldsymbol{\xi}_{1:T})$$

Smoothed edge marginals

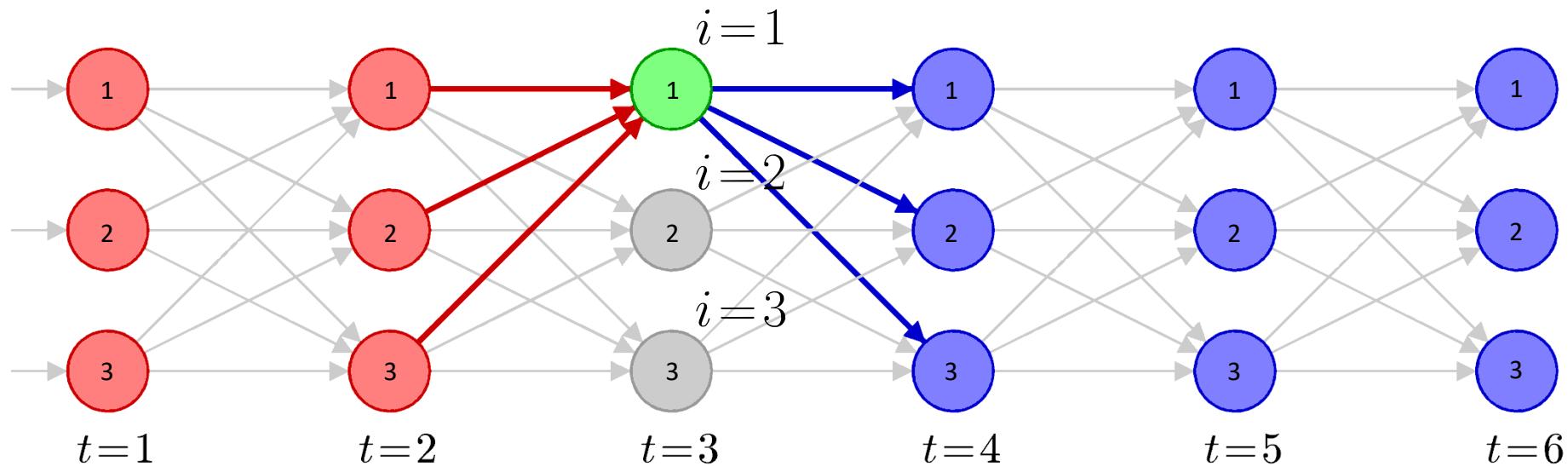
$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j | \boldsymbol{\xi}_{1:T})$$

Smoothed node marginals

$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i \mid \xi_{1:T})$$

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$$

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} \mid s_t=i)$$



$$\gamma_{t,i}^{\text{HMM}} = \frac{\alpha_{t,i}^{\text{HMM}} \beta_{t,i}^{\text{HMM}}}{\sum_{k=1}^K \alpha_{t,k}^{\text{HMM}} \beta_{t,k}^{\text{HMM}}} = \frac{\alpha_{t,i}^{\text{HMM}} \beta_{t,i}^{\text{HMM}}}{\mathcal{P}(\xi)}$$

Useful intermediary variables in HMM

Forward variable

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$$

Backward variable

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} \mid s_t=i)$$

Smoothed node marginals

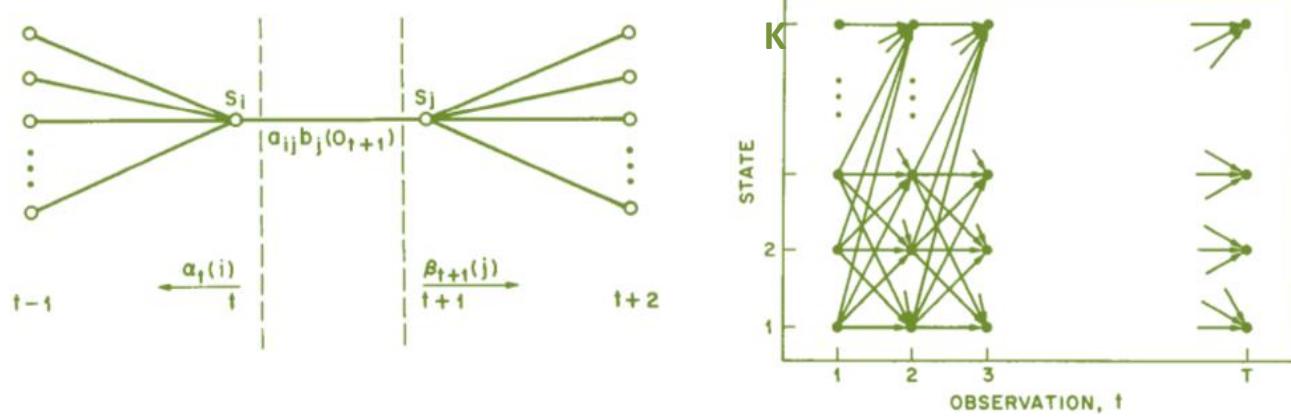
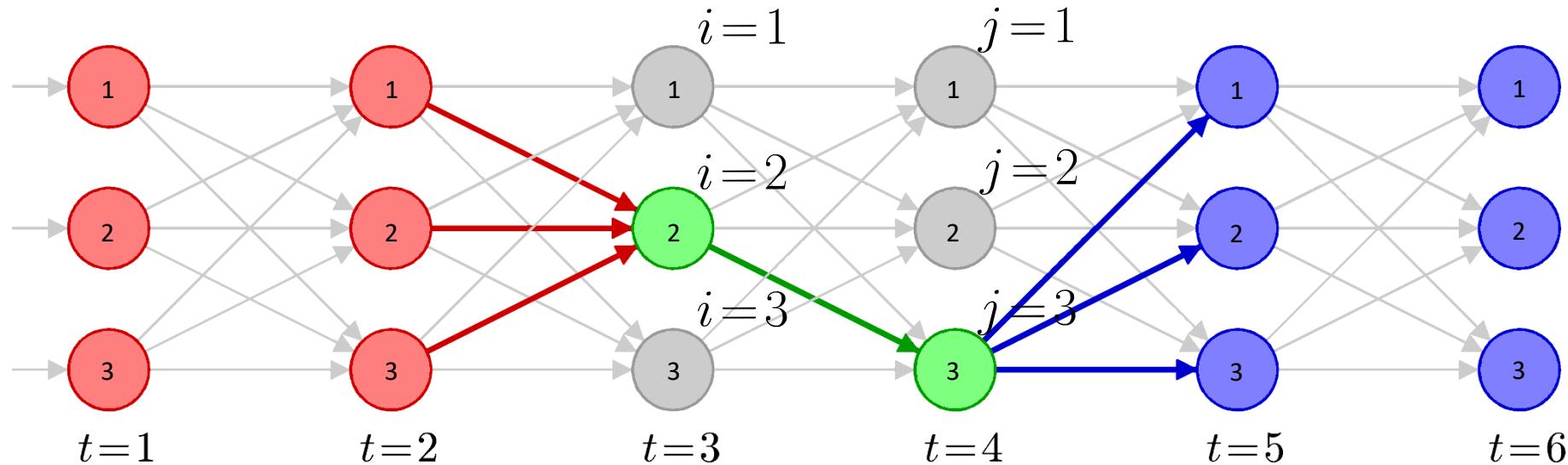
$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i \mid \xi_{1:T})$$

Smoothed edge marginals

$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j \mid \xi_{1:T})$$

Smoothed edge marginals

$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j | \xi_{1:T})$$



$$\begin{aligned}\zeta_{t,i,j}^{\text{HMM}} &= \frac{\alpha_{t,i}^{\text{HMM}} a_{i,j} \mathcal{N}(\xi_{t+1} | \mu_j, \Sigma_j) \beta_{t+1,j}^{\text{HMM}}}{\sum_{k=1}^K \sum_{l=1}^K \alpha_{t,k}^{\text{HMM}} a_{k,l} \mathcal{N}(\xi_{t+1} | \mu_l, \Sigma_l) \beta_{t+1,l}^{\text{HMM}}} \\ &= \frac{\alpha_{t,i}^{\text{HMM}} a_{i,j} \mathcal{N}(\xi_{t+1} | \mu_j, \Sigma_j) \beta_{t+1,j}^{\text{HMM}}}{\mathcal{P}(\xi)}\end{aligned}$$

EM for HMM (Baum-Welch algorithm)

E-step: compute $\gamma_{m,t,i}^{\text{HMM}}$ and $\zeta_{m,t,i,j}^{\text{HMM}}$

M-step:

$$\Pi_i \leftarrow \frac{\sum_{m=1}^M \gamma_{m,1,i}^{\text{HMM}}}{M} = \frac{\text{Total number of times in } i \text{ at time step 1}}{\text{Total number of trajectories}}$$

$$a_{i,j} \leftarrow \frac{\sum_{m=1}^M \sum_{t=1}^{T_m-1} \zeta_{m,t,i,j}^{\text{HMM}}}{\sum_{m=1}^M \sum_{t=1}^{T_m-1} \gamma_{m,t,i}^{\text{HMM}}} = \frac{\text{Total number of transitions from } i \text{ to } j}{\text{Total number of times in } i \text{ (and transit to anything else)}}$$

$$\boldsymbol{\mu}_i \leftarrow \frac{\sum_{m=1}^M \sum_{t=1}^{T_m} \gamma_{m,t,i}^{\text{HMM}} \boldsymbol{\xi}_{m,t}}{\sum_{m=1}^M \sum_{t=1}^{T_m} \gamma_{m,t,i}^{\text{HMM}}} \quad \text{result similar to GMM}$$

$$\boldsymbol{\Sigma}_i \leftarrow \frac{\sum_{m=1}^M \sum_{t=1}^{T_m} \gamma_{m,t,i}^{\text{HMM}} (\boldsymbol{\xi}_{m,t} - \boldsymbol{\mu}_i)(\boldsymbol{\xi}_{m,t} - \boldsymbol{\mu}_i)^\top}{\sum_{m=1}^M \sum_{t=1}^{T_m} \gamma_{m,t,i}^{\text{HMM}}}$$

K Gaussians

M sequences

T_m points per sequences

The update rules can be interpreted as normalized counts, with several types of weighted averages required in the computation.

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t = i, \boldsymbol{\xi}_{1:t})$$

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\boldsymbol{\xi}_{t+1:T} | s_t = i)$$

$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t = i | \boldsymbol{\xi}_{1:T})$$

$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t = i, s_{t+1} = j | \boldsymbol{\xi}_{1:T})$$

Numerical underflow issue in HMM

For long sequences, the forward and backward variables can quickly get very low, likely exceeding the precision range of the computer.

A simple scaling procedure is to multiply $\alpha_{t,i}^{\text{HMM}}$ by a factor independent of i , and divide $\beta_{t,i}^{\text{HMM}}$ by the same factor so that they are cancelled in the forward-backward computation.

The computation can be kept within reasonable bounds by setting the scaling factor

$$c_t = \frac{1}{\sum_{i=1}^K \alpha_{t,i}^{\text{HMM}}}$$

Numerical underflow issue in HMM

$$c_t = \frac{1}{\sum_{i=1}^K \alpha_{t,i}^{\text{HMM}}}$$

This issue is sometimes not covered in textbooks, although it remains very important for practical implementation of HMM!

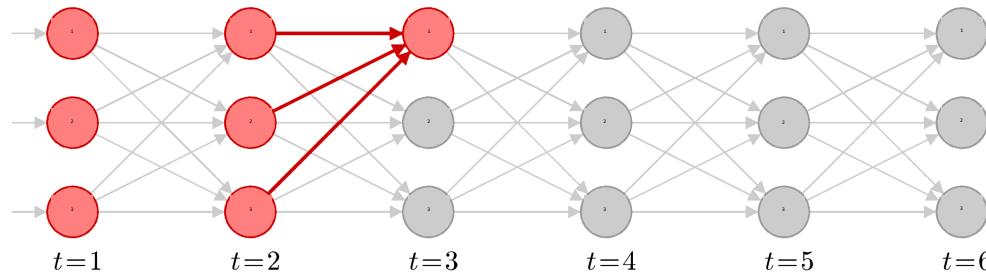
We have by induction

$$\hat{\alpha}_{t,i}^{\text{HMM}} = \left(\prod_{s=1}^t c_s \right) \alpha_{t,i}^{\text{HMM}}, \quad \hat{\beta}_{t,i}^{\text{HMM}} = \left(\prod_{s=t}^T c_s \right) \beta_{t,i}^{\text{HMM}}$$

With this, the numerator and denominator will cancel out when used in the re-estimation formulas. For example

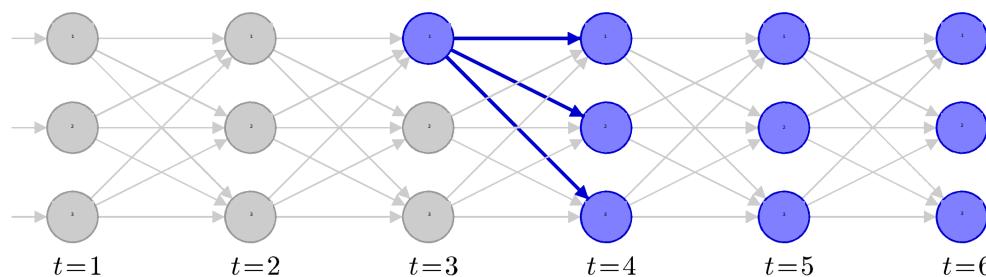
$$\gamma_{t,i}^{\text{HMM}} = \frac{\hat{\alpha}_{t,i}^{\text{HMM}} \hat{\beta}_{t,i}^{\text{HMM}}}{\sum_{k=1}^K \hat{\alpha}_{t,k}^{\text{HMM}} \hat{\beta}_{t,k}^{\text{HMM}}} = \frac{\left(\prod_{s=1}^t c_s \right) \left(\prod_{s=t}^T c_s \right) \alpha_{t,i}^{\text{HMM}} \beta_{t,i}^{\text{HMM}}}{\left(\prod_{s=1}^t c_s \right) \left(\prod_{s=t}^T c_s \right) \sum_{k=1}^K \alpha_{t,k}^{\text{HMM}} \beta_{t,k}^{\text{HMM}}} = \frac{\alpha_{t,i}^{\text{HMM}} \beta_{t,i}^{\text{HMM}}}{\sum_{k=1}^K \alpha_{t,k}^{\text{HMM}} \beta_{t,k}^{\text{HMM}}}$$

Why did we introduce these intermediary variables in HMM?



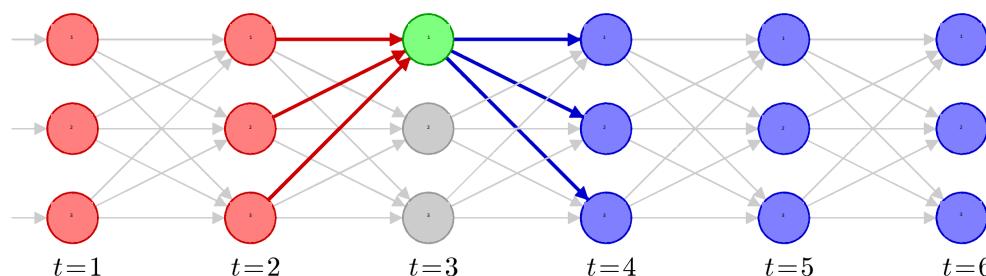
Forward variable

$$\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$$



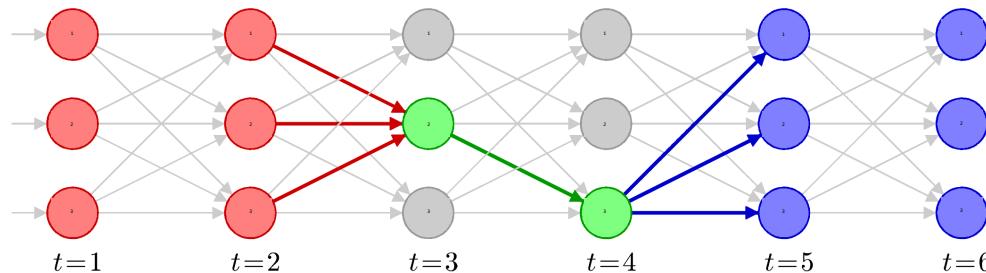
Backward variable

$$\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} \mid s_t=i)$$



Smoothed node marginals

$$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i \mid \xi_{1:T})$$



Smoothed edge marginals

$$\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j \mid \xi_{1:T})$$

Why did we introduce these intermediary variables in HMM?

How to estimate the parameters of an HMM?

→ Maximum of expected complete data log-likelihood $\mathcal{Q}(\Theta, \Theta^{\text{old}})$

How to compute $\frac{\partial \mathcal{Q}}{\partial \Pi_i} = 0$, $\frac{\partial \mathcal{Q}}{\partial a_{i,j}} = 0$, $\frac{\partial \mathcal{Q}}{\partial \mu_i} = 0$ and $\frac{\partial \mathcal{Q}}{\partial \Sigma_i} = 0$?

→ Requires to compute $\zeta_{t,i,j}^{\text{HMM}} = \mathcal{P}(s_t=i, s_{t+1}=j | \xi_{1:T})$

→ Requires to compute $\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i | \xi_{1:T})$

How to compute $\zeta_{t,i,j}^{\text{HMM}}$ and $\gamma_{t,i}^{\text{HMM}}$?

→ Requires to compute $\alpha_{t,i}^{\text{HMM}} = \mathcal{P}(s_t=i, \xi_{1:t})$

→ Requires to compute $\beta_{t,i}^{\text{HMM}} = \mathcal{P}(\xi_{t+1:T} | s_t=i)$

$$\max_{\Theta} \mathcal{Q}(\Theta, \Theta^{\text{old}})$$

$$\zeta_{t,i,j}^{\text{HMM}} \quad \gamma_{t,i}^{\text{HMM}}$$

$$\alpha_{t,i}^{\text{HMM}}$$

$$\beta_{t,i}^{\text{HMM}}$$

Viterbi decoding (MAP vs MPE estimates)

Maximum a posteriori

Most probable explanation

Python notebook:
`demo_HMM.ipynb`

Matlab code:
`demo_HMM_Viterbi01.m`

Viterbi decoding (MAP vs MPE estimates)

Maximum a posteriori

Most probable explanation

The (jointly) most probable sequence of states \hat{s}^{MAP} is not necessarily the same as the sequence of (marginally) most probable states \hat{s}^{MPE}

$$\hat{s}^{\text{MAP}} = \arg \max_{\{s_1, s_2, \dots, s_T\}} \mathcal{P}(s | \xi)$$

$$\hat{s}^{\text{MPE}} = \left\{ \arg \max_{s_1} \overbrace{\mathcal{P}(s_1 | \xi)}^{\gamma_1^{\text{HMM}}}, \arg \max_{s_2} \overbrace{\mathcal{P}(s_2 | \xi)}^{\gamma_2^{\text{HMM}}}, \dots, \arg \max_{s_T} \overbrace{\mathcal{P}(s_T | \xi)}^{\gamma_T^{\text{HMM}}} \right\}$$

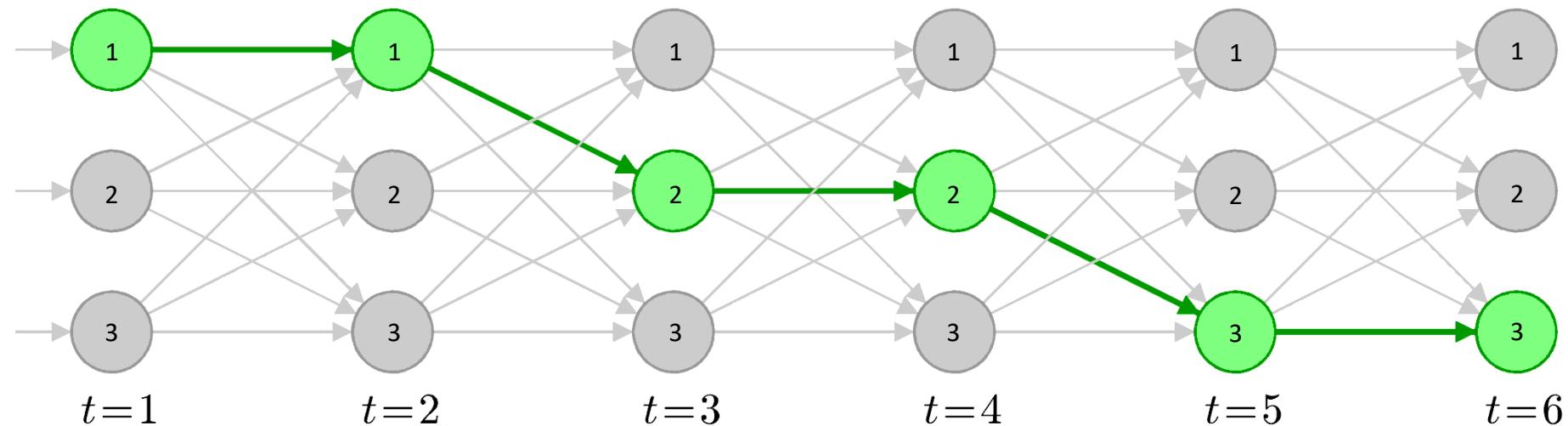
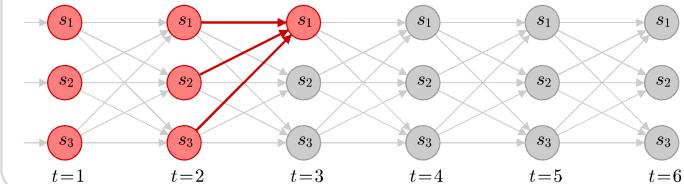
$\gamma_{t,i}^{\text{HMM}} = \mathcal{P}(s_t = i | \xi_{1:T})$

\hat{s}^{MAP} can be computed with the **Viterbi algorithm**, employing the max operator in a forward pass, followed by a backward pass using a **fast traceback procedure** to recover the most probable path.

Viterbi decoding

Forward algorithm:

$$\alpha_{t,i}^{\text{HMM}} = \left(\sum_{j=1}^K \alpha_{t-1,j}^{\text{HMM}} a_{j,i} \right) \mathcal{N}(\xi_t | \mu_i, \Sigma_i)$$



Forward recursion:

$$\begin{cases} \delta_{t,i} = \max_j (\delta_{t-1,j} a_{j,i}) \mathcal{N}(\xi_t | \mu_i, \Sigma_i) & \leftarrow \text{probability of ending in state } i \text{ at time step } t \text{ by taking the most probable path, with } \delta_{1,i} = \Pi_i \mathcal{N}(\xi_1 | \mu_i, \Sigma_i) \\ \Psi_{t,i} = \arg \max_j (\delta_{t-1,j} a_{j,i}) & \leftarrow \text{indices that keep track of the states } j \text{ that maximized } \delta_{t,i}. \end{cases}$$

Backtracking: $\hat{s}_t^{\text{MAP}} = \Psi_{t+1, \hat{s}_{t+1}^{\text{MAP}}}$

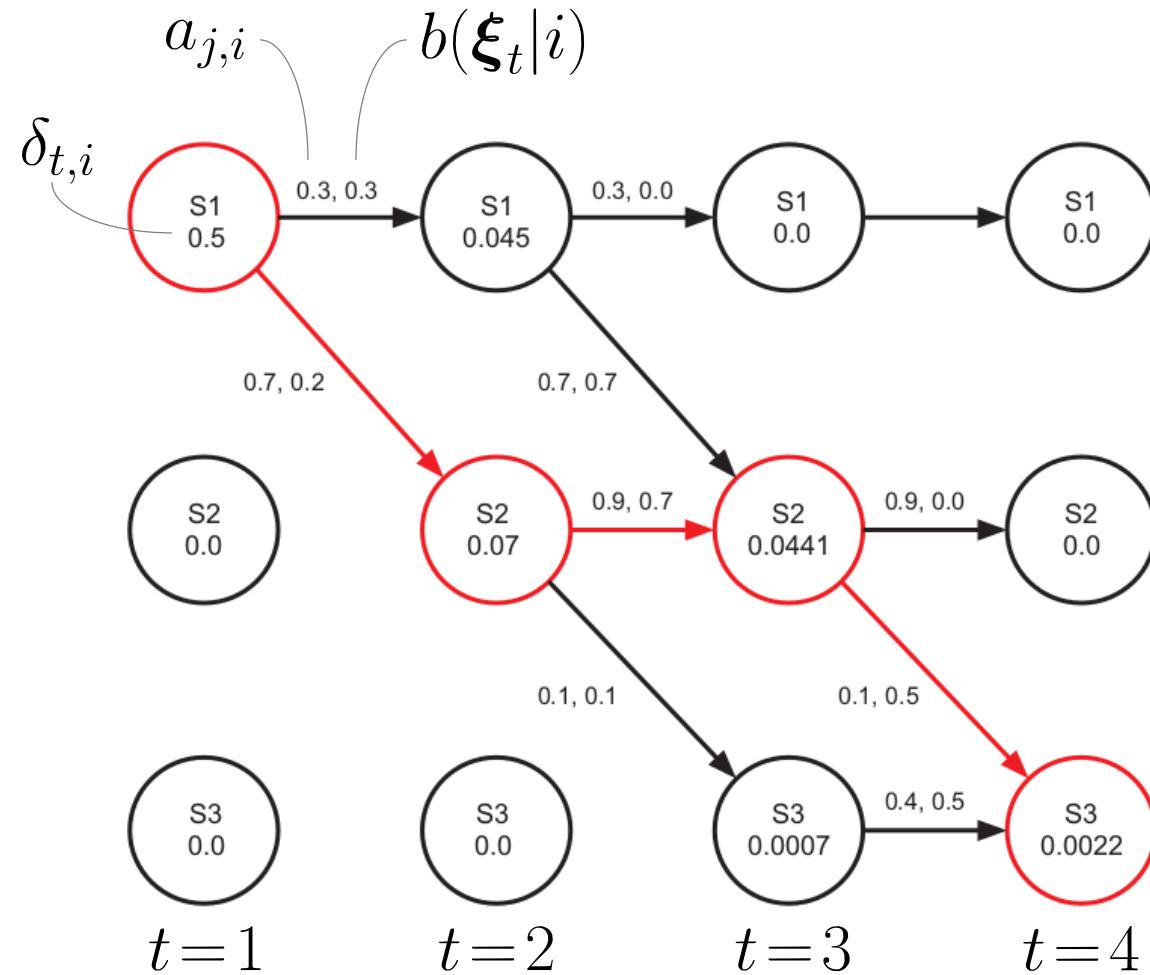
Viterbi decoding - Example

$$\delta_{t,i} = \max_j (\delta_{t-1,j} a_{j,i}) b(\xi_t | i)$$

with $\delta_{1,i} = \Pi_i \mathcal{N}(\xi_1 | \mu_i, \Sigma_i)$

$$\Pi = [1, 0, 0]$$

Dictionary	$S1$	0.3	0.7	$S2$	0.9	0.1	$S3$	0.4
	$C1$	0.5	0	0				
	$C2$	0.2	0	0				
	$C3$	0.3	0.2	0				
	$C4$	0	0.7	0.1				
	$C5$	0	0.1	0				
	$C6$	0	0	0.5				
	$C7$	0	0	0.4				



Observation: $\xi = \{C1, C3, C4, C6\}$

Numerical underflow issue in Viterbi

Similarly to the forward-backward variables in HMM, we have to take care about potential numerical underflow when implementing Viterbi decoding.

A simple way is to normalize $\delta_{t,i}$ at each time step t with

$$c_t = \frac{1}{\sum_{i=1}^K \delta_{t,i}}$$

similarly as in the computation of the forward-backward variables.

Such scaling will not affect the maximum.

Numerical underflow issue in Viterbi

Alternatively, we can work in the log domain. We then have

$$\begin{aligned}\log \delta_{t,i} &= \max_{\boldsymbol{s}_{1:t-1}} \log \mathcal{P}(\boldsymbol{s}_{1:t-1}, s_t = i \mid \boldsymbol{\xi}_{1:t}) \\ &= \max_j (\log \delta_{t-1,j} + \log a_{i,j}) + \log \mathcal{N}(\boldsymbol{\xi}_t \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\end{aligned}$$

With high dimensional Gaussians as emission distributions, the Viterbi computation with log can result in a **significant speedup**, since computing $\log \mathcal{P}(\boldsymbol{\xi}_t | s_t)$ can be much faster than computing $\mathcal{P}(\boldsymbol{\xi}_t | s_t)$.

When training HMMs, the Viterbi algorithm can be also used (instead of the forward-backward variables) in the E step of the EM procedure.

Hidden semi-Markov model (HSMM)

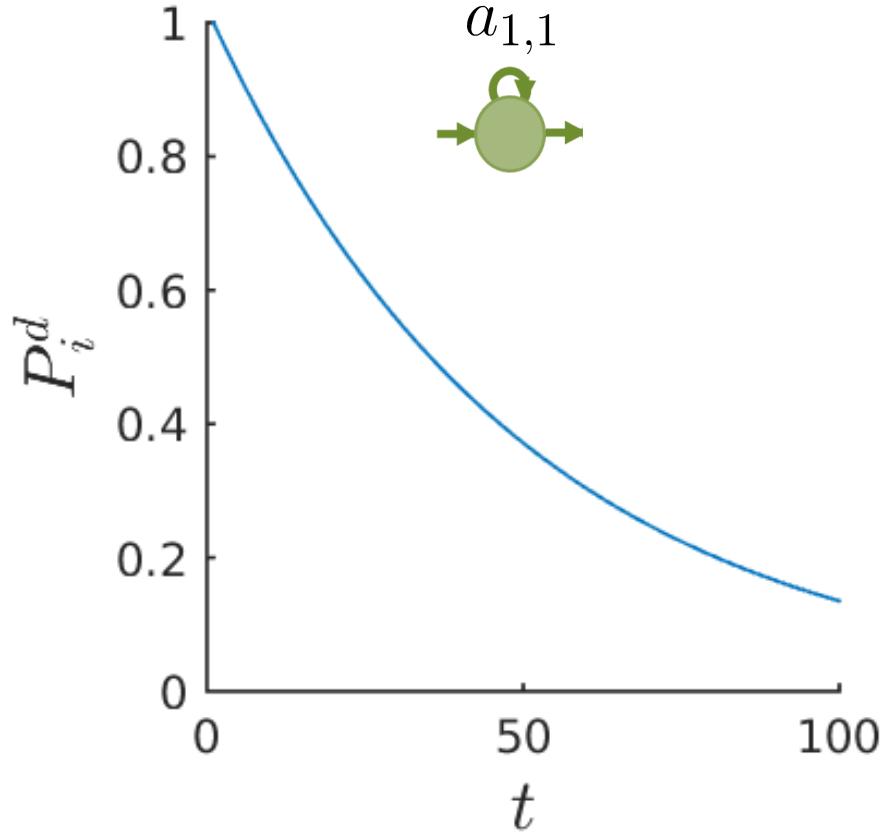
Python notebook:
`demo_HSMM.ipynb`

Matlab code:
`demo_HSMM01.m`

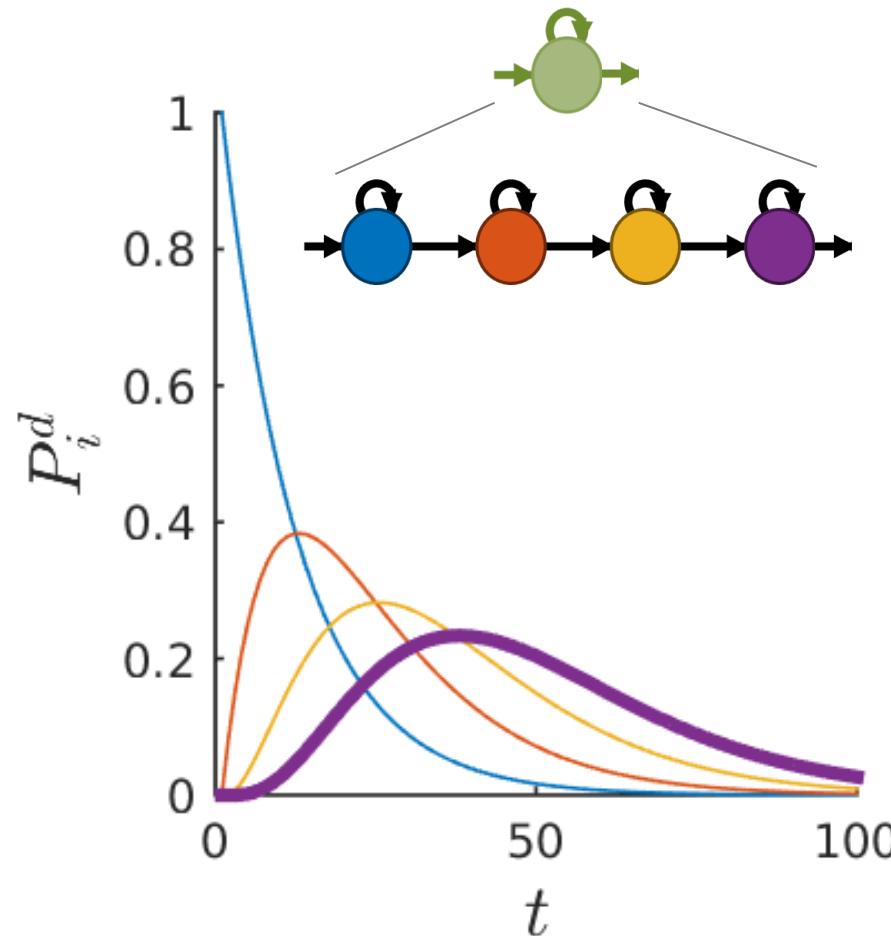
State duration probability in standard HMM

The state duration follows a geometric distribution

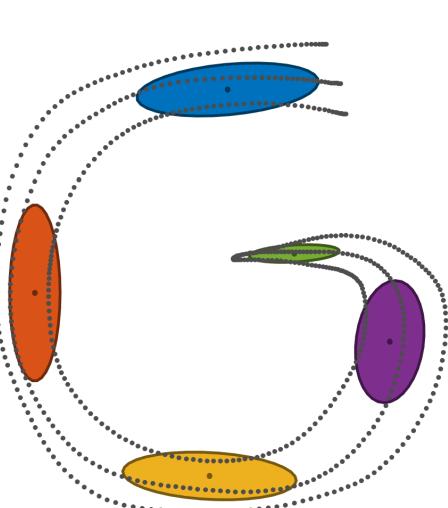
$$\mathcal{P}(d) = a_{i,i}^{d-1}(1 - a_{i,i})$$



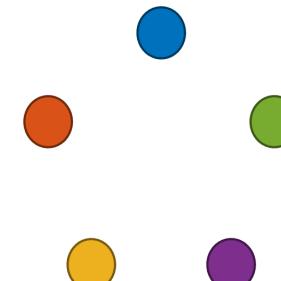
Naive approach: By artificially duplicating the number of states while keeping the same emission distribution, other state duration distributions can be modeled



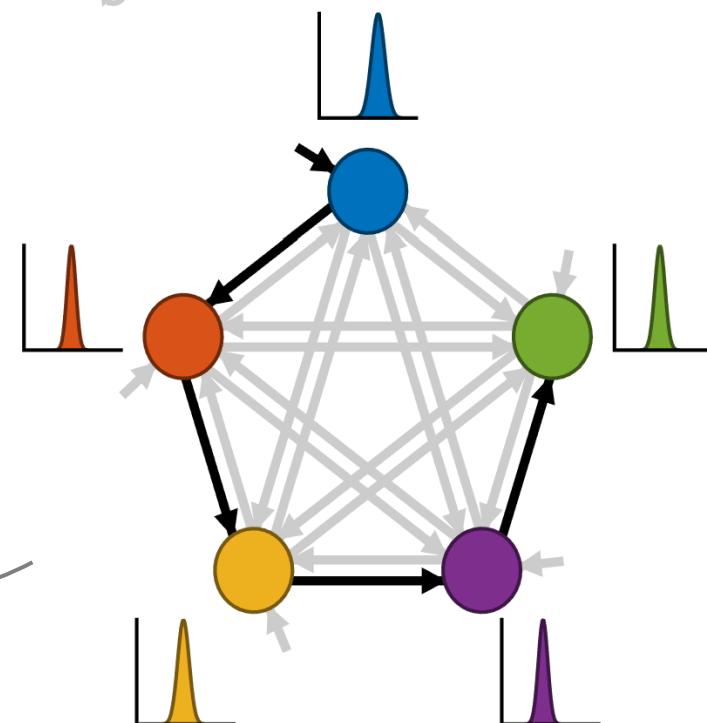
Hidden semi-Markov model (HSMM)



GMM



HMM

HSMM

Another approach is to provide an explicit model of the state duration instead of relying on self-transition probabilities

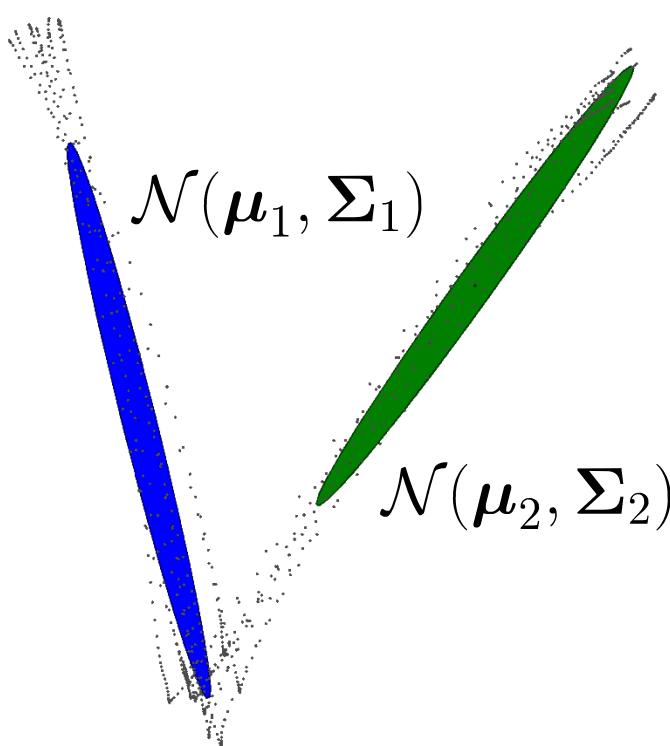
Hidden semi-Markov model (HSMM)

$$\Theta^{\text{GMM}} = \{\pi_i, \mu_i, \Sigma_i\}_{i=1}^K$$

$$\Theta^{\text{HMM}} = \{\{a_{i,j}\}_{j=1}^K, \Pi_i, \mu_i, \Sigma_i\}_{i=1}^K$$

$$\Theta^{\text{HSMM}} = \{\{a_{i,j}\}_{j=1, j \neq i}^K, \Pi_i, \mu_i^{\mathcal{D}}, \Sigma_i^{\mathcal{D}}, \mu_i, \Sigma_i\}_{i=1}^K$$

Parametric duration
distribution

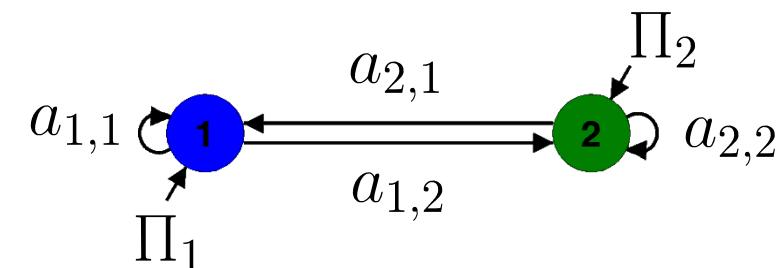


GMM

1

2

HMM

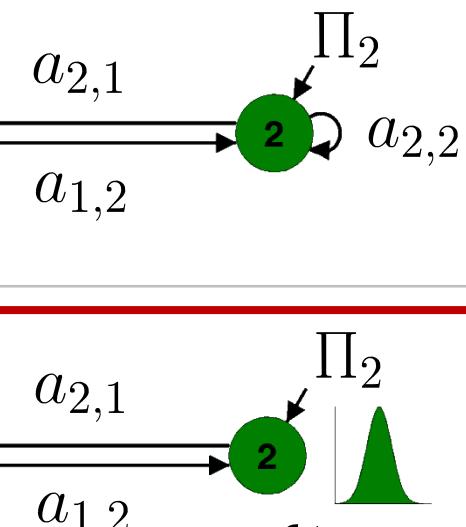


HSMM

$\mathcal{N}(\mu_1^{\mathcal{D}}, \Sigma_1^{\mathcal{D}})$

Π_1

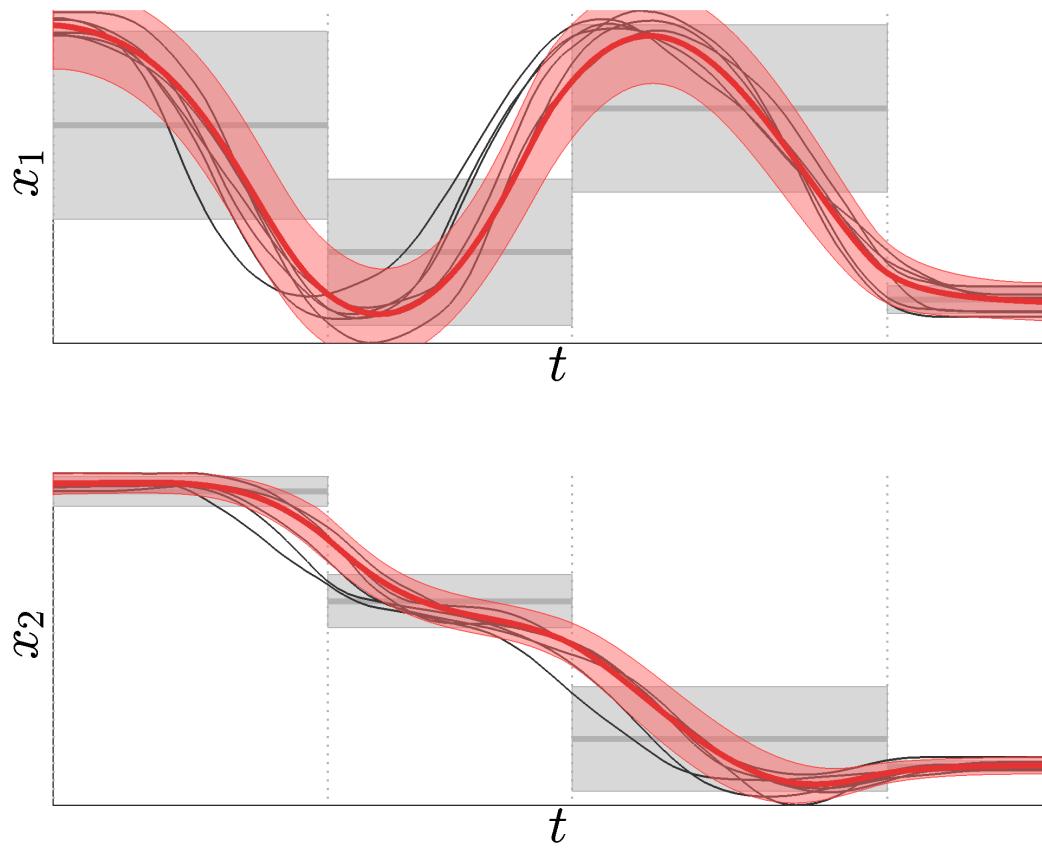
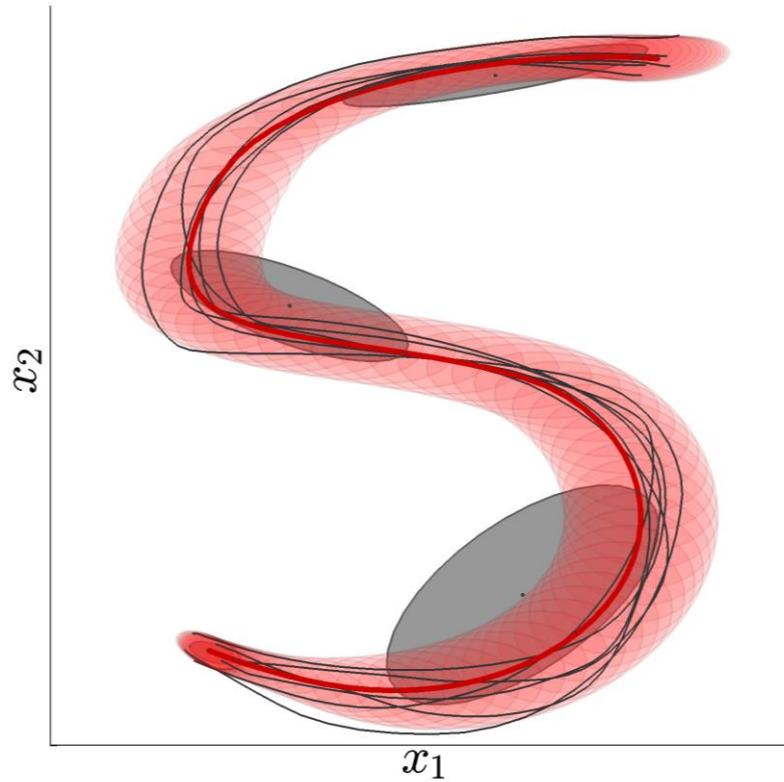
$\mathcal{N}(\mu_2^{\mathcal{D}}, \Sigma_2^{\mathcal{D}})$



HMM with dynamic features (Trajectory-HMM)

Matlab code:
`demo_trajHSMM01.m`

HMM with dynamic features



HMM with dynamic features

For the encoding of movements, velocity and acceleration can be used as dynamic features. By considering an Euler approximation, the velocity is computed as

$$\dot{\mathbf{x}}_t = \frac{\mathbf{x}_{t+1} - \mathbf{x}_t}{\Delta t}$$

where \mathbf{x}_t is a multivariate position vector.

The acceleration is similarly computed as

$$\ddot{\mathbf{x}}_t = \frac{\dot{\mathbf{x}}_{t+1} - \dot{\mathbf{x}}_t}{\Delta t} = \frac{\mathbf{x}_{t+2} - 2\mathbf{x}_{t+1} + \mathbf{x}_t}{\Delta t^2}$$

HMM with dynamic features

$$\dot{\boldsymbol{x}}_t = \frac{\boldsymbol{x}_{t+1} - \boldsymbol{x}_t}{\Delta t}, \quad \ddot{\boldsymbol{x}}_t = \frac{\boldsymbol{x}_{t+2} - 2\boldsymbol{x}_{t+1} + \boldsymbol{x}_t}{\Delta t^2}$$

A vector ζ_t will be used to represent the concatenated position, velocity and acceleration vectors at time step t

$$\zeta_t = \begin{bmatrix} \boldsymbol{x}_t \\ \dot{\boldsymbol{x}}_t \\ \ddot{\boldsymbol{x}}_t \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ -\frac{1}{\Delta t} \mathbf{I} & \frac{1}{\Delta t} \mathbf{I} & \mathbf{0} \\ \frac{1}{\Delta t^2} \mathbf{I} & -\frac{2}{\Delta t^2} \mathbf{I} & \frac{1}{\Delta t^2} \mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_t \\ \boldsymbol{x}_{t+1} \\ \boldsymbol{x}_{t+2} \end{bmatrix}$$

Here, the number of derivatives will be set up to acceleration ($C=3$), but the same approach can be applied to a different number of derivatives.

A GMM/HMM/HSMM with centers $\{\boldsymbol{\mu}_i\}_{i=1}^K$ and covariance matrices $\{\boldsymbol{\Sigma}_i\}_{i=1}^K$ is first fit to the dataset $[\zeta_1, \zeta_2, \dots, \zeta_T]$.

HMM with dynamic features

$$\zeta_t = \begin{bmatrix} \mathbf{x}_t \\ \dot{\mathbf{x}}_t \\ \ddot{\mathbf{x}}_t \end{bmatrix}$$

ζ and \mathbf{x} are defined as large vectors concatenating ζ_t and \mathbf{x}_t for all time steps

$$\zeta = \begin{bmatrix} \zeta_1 \\ \zeta_2 \\ \vdots \\ \zeta_T \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}$$

Similarly to the matrix operator defined in the previous slide for a single time step, a large sparse matrix Φ can be defined so that

$$\zeta = \Phi \mathbf{x}$$

HMM with dynamic features

D dimensions
 C derivatives
 T time steps

$$\begin{bmatrix} \vdots \\ \boldsymbol{x}_t \\ \dot{\boldsymbol{x}}_t \\ \ddot{\boldsymbol{x}}_t \\ \boldsymbol{x}_{t+1} \\ \dot{\boldsymbol{x}}_{t+1} \\ \ddot{\boldsymbol{x}}_{t+1} \\ \vdots \end{bmatrix} = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots & \dots \\ \dots & -\frac{1}{\Delta t} \mathbf{I} & \frac{1}{\Delta t} \mathbf{I} & \mathbf{0} & \dots & \dots \\ \dots & \frac{1}{\Delta t^2} \mathbf{I} & -\frac{1}{\Delta t^2} \mathbf{I} & \frac{1}{\Delta t^2} \mathbf{I} & \dots & \dots \\ \dots & \dots & \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & \dots & -\frac{1}{\Delta t} \mathbf{I} & \frac{1}{\Delta t} \mathbf{I} & \mathbf{0} & \dots \\ \dots & \dots & \frac{1}{\Delta t^2} \mathbf{I} & -\frac{1}{\Delta t^2} \mathbf{I} & \frac{1}{\Delta t^2} \mathbf{I} & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \boldsymbol{x}_t \\ \boldsymbol{x}_{t+1} \\ \boldsymbol{x}_{t+2} \\ \boldsymbol{x}_{t+3} \\ \vdots \end{bmatrix}$$

$\zeta \in \mathbb{R}^{DCT}$

(C=3 here)

$\Phi \in \mathbb{R}^{DCT \times DT}$

$\boldsymbol{x} \in \mathbb{R}^{DT}$

Large sparse matrix

HMM with dynamic features

By providing explicitly a sequence of states $\mathbf{s} = \{s_1, s_2, \dots, s_T\}$ of T time steps (e.g., retrieved by Viterbi or specified manually), with discrete states $s_t \in \{1, \dots, K\}$, the likelihood of a movement $\zeta = \Phi \mathbf{x}$ is given by

$$\mathcal{P}(\zeta | \mathbf{s}) = \prod_{t=1}^T \mathcal{N}(\zeta_t | \boldsymbol{\mu}_{s_t}, \Sigma_{s_t})$$

where $\boldsymbol{\mu}_{s_t}$ and Σ_{s_t} are the center and covariance of state s_t at time step t .

This product can be rewritten as

$$\mathcal{P}(\zeta | \mathbf{s}) = \mathcal{N}(\zeta | \boldsymbol{\mu}_s, \Sigma_s)$$

with $\boldsymbol{\mu}_s = \begin{bmatrix} \boldsymbol{\mu}_{s_1} \\ \boldsymbol{\mu}_{s_2} \\ \vdots \\ \boldsymbol{\mu}_{s_T} \end{bmatrix}$ and $\Sigma_s = \begin{bmatrix} \Sigma_{s_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma_{s_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma_{s_T} \end{bmatrix}$

HMM with dynamic features

For example, for a sequence of states $\mathbf{s} = \{1, 1, 2, 2, 3, 3, 3, 4\}$ with $K=4$ and $T=8$, we have

$$\mu_s = \begin{bmatrix} \mu_1 \\ \mu_1 \\ \mu_2 \\ \mu_2 \\ \mu_3 \\ \mu_3 \\ \mu_3 \\ \mu_4 \end{bmatrix} \quad \text{and} \quad \Sigma_s = \begin{bmatrix} \Sigma_1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \Sigma_1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \Sigma_2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \Sigma_2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \Sigma_3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \Sigma_3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \Sigma_3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \Sigma_4 \end{bmatrix}$$

$$\mu_s \in \mathbb{R}^{DCT}$$

$$\Sigma_s \in \mathbb{R}^{DCT \times DCT}$$

HMM with dynamic features

By using the relation $\zeta = \Phi x$, we want to retrieve a trajectory

$$\hat{x} = \arg \max_x \log \mathcal{P}(\Phi x | s)$$

$$\mathcal{N}(\Phi x | \mu_s, \Sigma_s) = (2\pi)^{-\frac{DCT}{2}} |\Sigma_s|^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\Phi x - \mu_s)^\top \Sigma_s^{-1} (\Phi x - \mu_s) \right)$$

Equating to zero the derivative of

$$\log \mathcal{P}(\Phi x | s) = -\frac{1}{2} (\Phi x - \mu_s)^\top \Sigma_s^{-1} (\Phi x - \mu_s) - \frac{1}{2} \log |\Sigma_s| - \frac{DCT}{2} \log(2\pi)$$

with respect to x yields

$$\Phi^\top \Sigma_s^{-1} (\Phi x - \mu_s) = 0$$

$$\iff \hat{x} = (\Phi^\top \Sigma_s^{-1} \Phi)^{-1} \Phi^\top \Sigma_s^{-1} \mu_s$$

Weighted least squares!

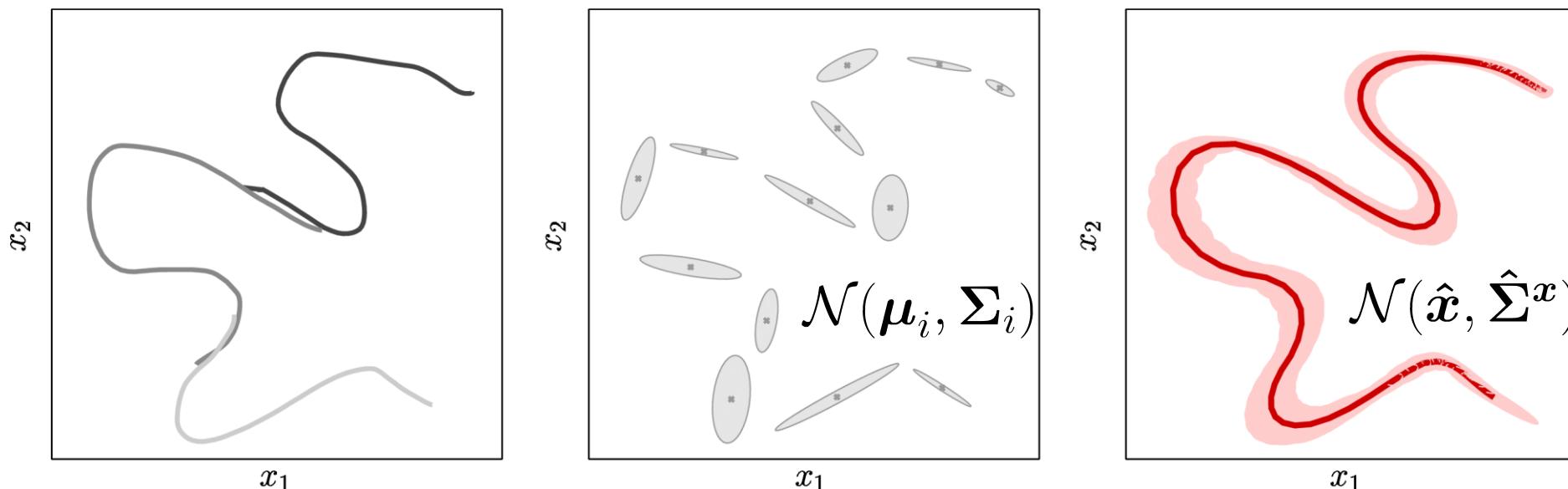
HMM with dynamic features

The residual error of this estimate is given by the covariance matrix

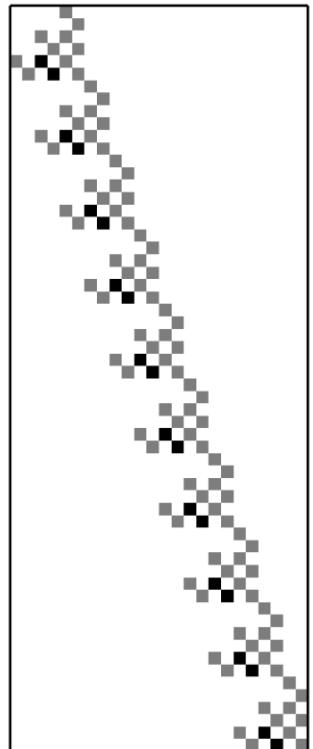
$$\hat{\Sigma}^x = \sigma (\Phi^\top \Sigma_s^{-1} \Phi)^{-1}$$

where σ is a scaling factor.

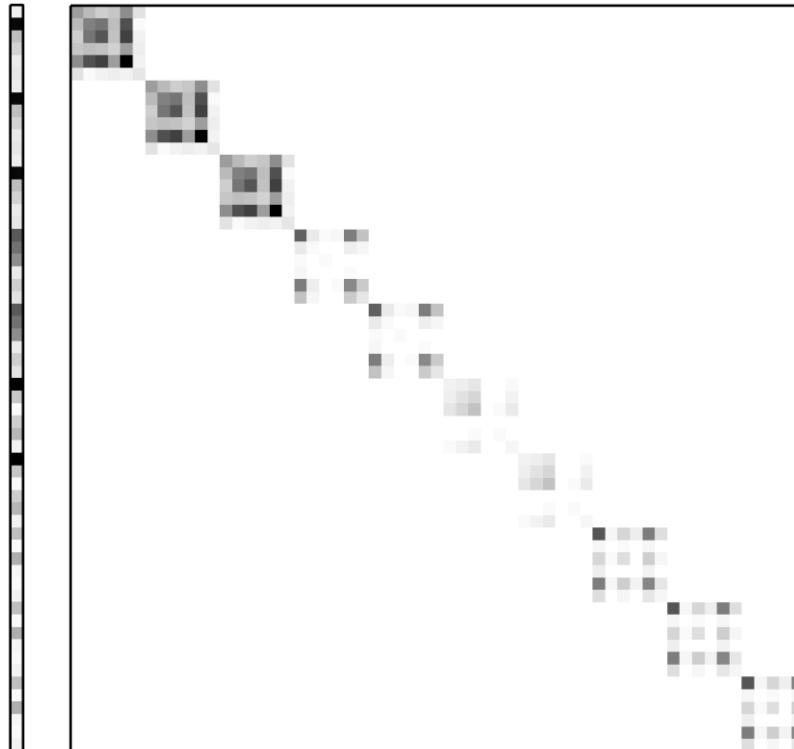
The resulting Gaussian $\mathcal{N}(\hat{x}, \hat{\Sigma}^x)$ forms a trajectory distribution, where $\hat{x} \in \mathbb{R}^{DT}$ is an average trajectory stored in a vector form.



HMM with dynamic features - Summary



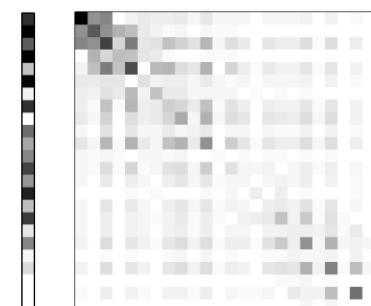
Φ



$\mathcal{N}(\mu_s, \Sigma_s)$

$$\hat{x} = (\Phi^\top \Sigma_s^{-1} \Phi)^{-1} \Phi^\top \Sigma_s^{-1} \mu_s$$

$$\hat{\Sigma}^x = \sigma (\Phi^\top \Sigma_s^{-1} \Phi)^{-1}$$



$\mathcal{N}(\hat{x}, \hat{\Sigma}^x)$

Reading material

Hidden Markov model (HMM)

L. R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. Proc. IEEE, 77:2:257–285, February 1989

Hidden semi-Markov model (HSMM)

S.-Z. Yu. Hidden semi-Markov models. Artificial Intelligence, 174:215–243, 2010

S. E. Levinson. Continuously variable duration hidden Markov models for automatic speech recognition. Computer Speech & Language, 1(1):29–45, 1986

HMM with dynamic features (Trajectory HMM)

S. Furui. Speaker-independent isolated word recognition using dynamic features of speech spectrum. IEEE Trans. on Acoustics, Speech, and Signal Processing, 34(1):52–59, 1986

H. Zen, K. Tokuda, and T. Kitamura. Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic feature vector sequences. Computer Speech and Language, 21(1):153–173, 2007