

“High accuracy and low computational requirement define a novel heuristic method for selecting the optimal order in hidden semi-Markov models”

Selection of the number of states in Hidden semi-Markov models

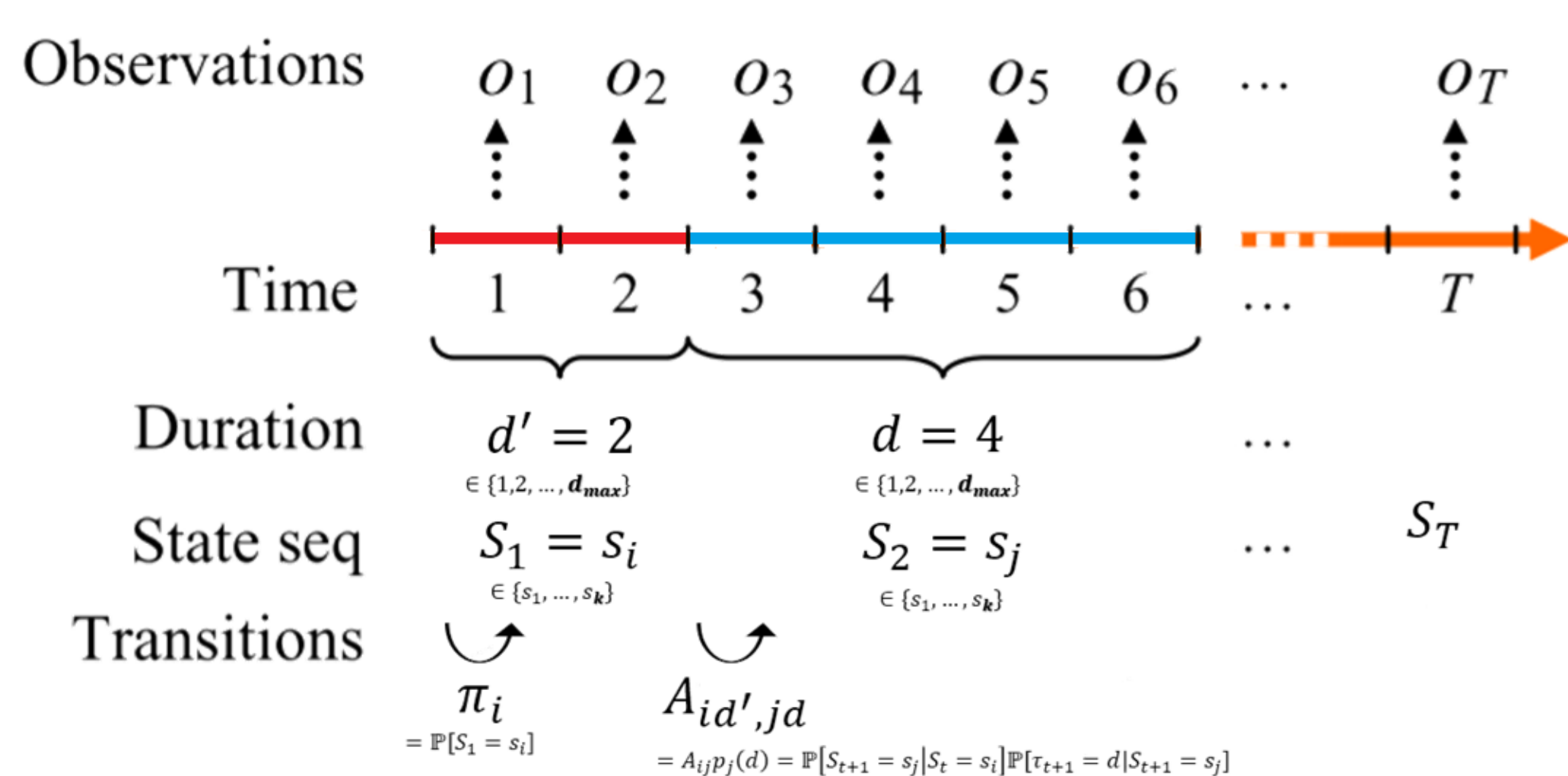
RENTE, Filipa¹; MARINHO, Zita²; FIGUEIREDO, Mário³

¹ Instituto Superior Técnico (IST), Universidade de Lisboa (ULisboa), Portugal

² Instituto de Sistemas e Robótica, IST & Priberam Labs, Lisboa, Portugal

³ Instituto de Telecomunicações, IST, ULisboa, Portugal

HIDDEN SEMI-MARKOV MODEL (HSMM)



CRITERIA ADAPTATION FOR HSMMs

BIC criterion for HSMMs

$$\text{BIC}_{\text{HSMM}}(k) = \text{LL} - \frac{k^2 + k}{2} \log(n)$$

k : number of states.

n : total number of samples.

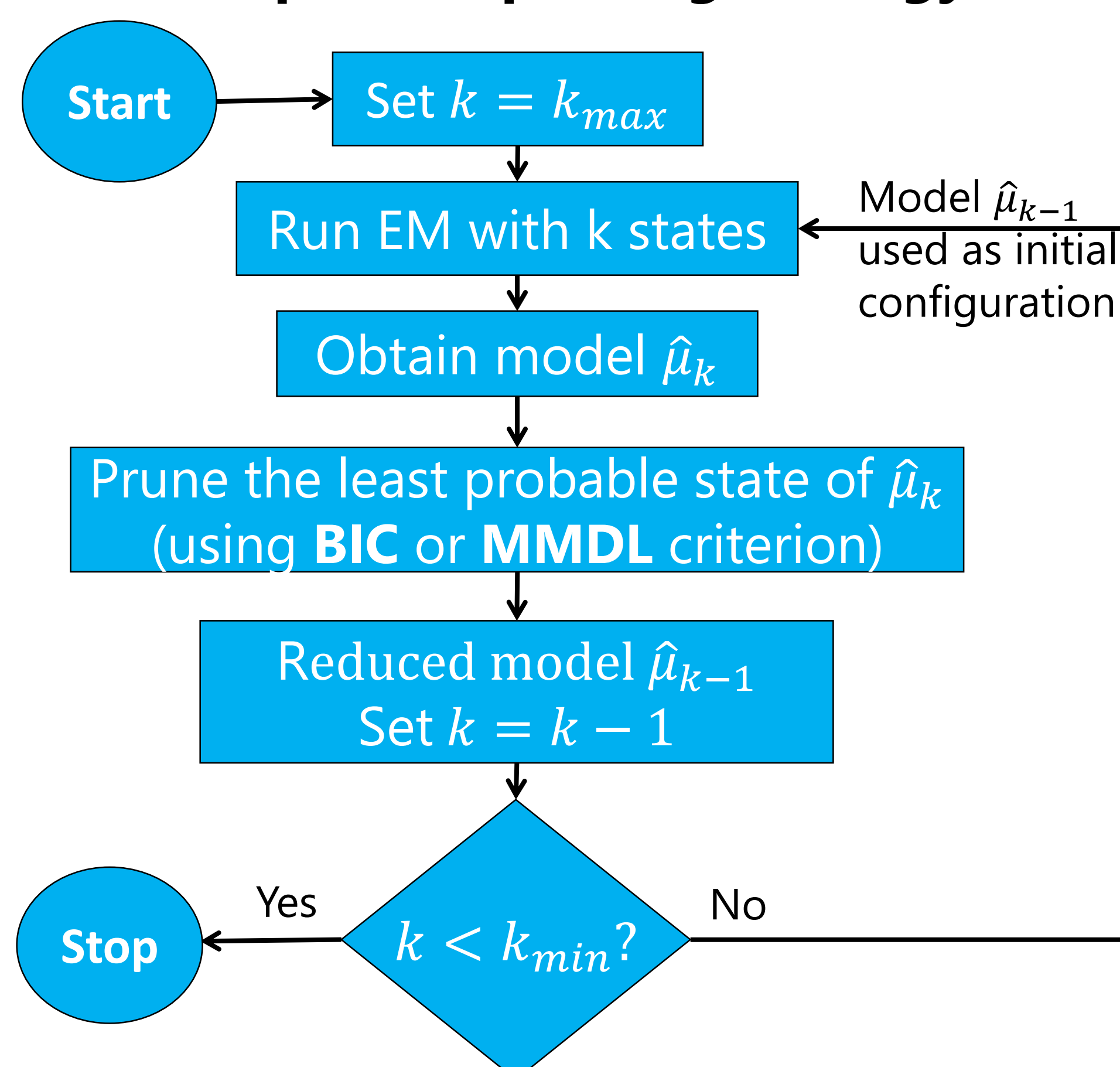
LL : log-likelihood of the observations.

MMDL criterion for HSMMs

$$\text{MMDL}_{\text{HSMM}}(k) = \text{LL} - \frac{k^2 - k}{2} \log(n) - \sum_{m=1}^k \log(np_{\infty}(m))$$

$p_{\infty}(m)$: stationary probability distribution for HSMMs.

Sequential pruning strategy



$\hat{\mu}_k = (s, A, \pi, B, D)$: predicted model parameters.

$s = \{s_1, \dots, s_k\}$: set of k states.

$A \in \mathbb{R}^{k \times k}$: transition matrix.

$\pi \in \mathbb{R}^k$: initial state probability distribution.

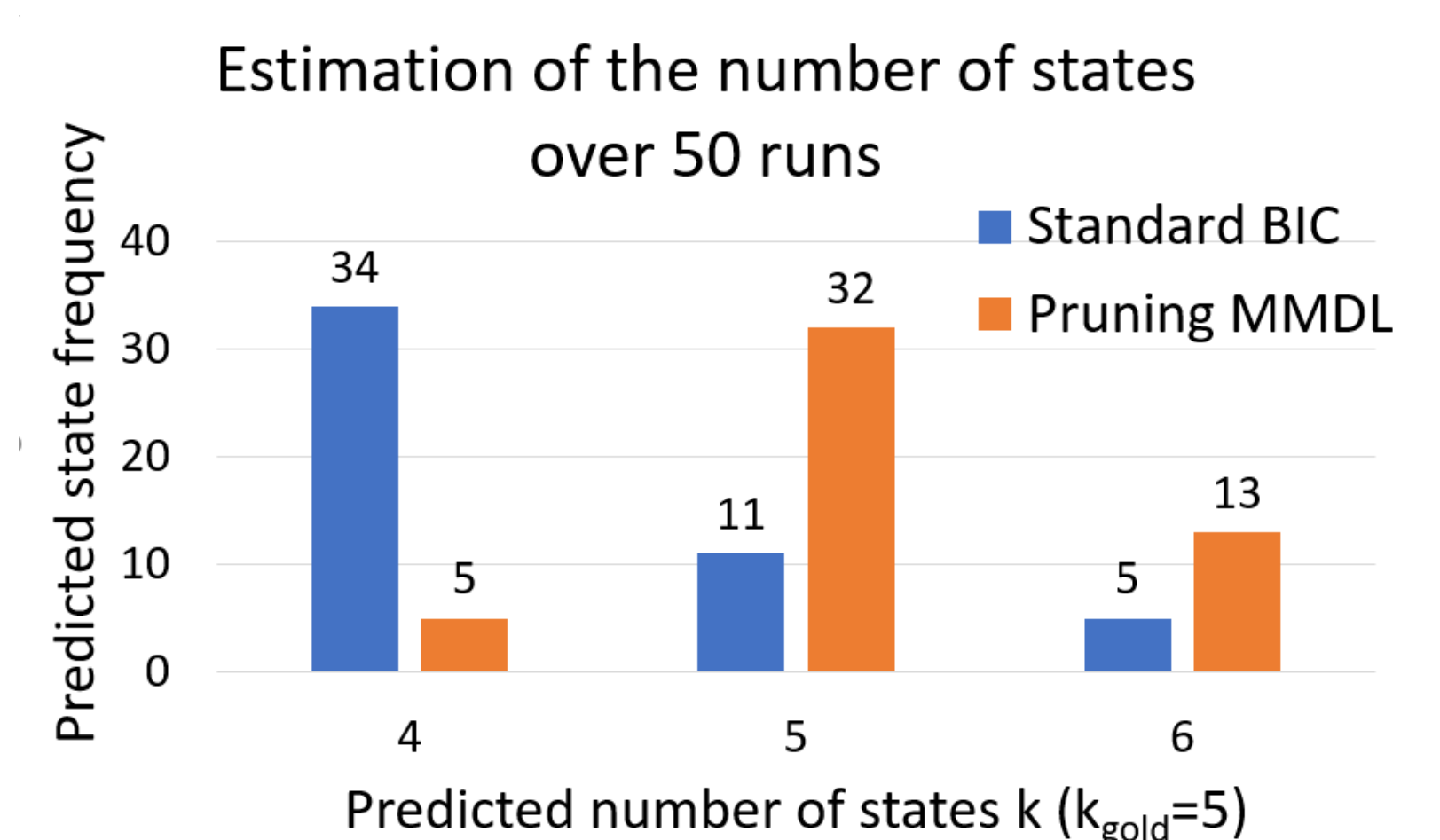
$B \in \mathbb{R}^{k \times m}$: emission matrix, for m different observation symbols.

$D \in \mathbb{R}^{k \times d_{\max}}$: duration probability matrix.

RESULTS with synthetic data

Pruning MMDL vs. BIC

Method	N ^o states	N ^o iter.
BIC	4.4 ± 0.7	38.6 ± 34.2
MMDL	5.2 ± 0.6	5.8 ± 2.9



Accuracy in the estimation of the number of states ($N^{\circ}states$ field; true value is 5);

Computational requirement reflected in the number of required iterations in the EM algorithm ($N^{\circ}iter$ field).

DISCUSSION

Comparison between the pruning MMDL for HSMMs and the standard BIC criterion:

Higher accuracy in the selection of the optimal number of states;

Less demanding **computational requirement**;

The sequential pruning strategy guarantees a **lower sensitivity** to the initialization of the EM algorithm.

FUTURE WORK

Compare with other standard methods;
Demonstrate the effectiveness in a more substantial application study using real data;

Design feature conditioned state transitions or add latent variables to the HSMM model so that external information can be captured.

INTRODUCTION

- HSMMs are powerful probabilistic models used in several fields, e.g., speech, health and genetics.
- More expressive extension of hidden Markov models (HMMs) with explicit state duration.
- Order selection problem: find the optimal hyperparameters k (n° states of the model) and d_{\max} (maximum allowed state duration).

SOLUTION

- Single HSMM selection criterion for both the number of states and the maximum allowed state duration.

METHODS

- Empiric estimation of d_{\max} :
 - Find θ_j (parameter of the state duration distribution) that corresponds to the highest duration amongst all states $s_j - \theta_j^*$;
 - d_{\max} is the (first) integer duration for which the cumulative state duration probability, with parameter θ_j^* , is smaller than $\varepsilon = 0,01$.

- Estimation of k :

Sequential pruning strategy
+
Mixture minimum description length (MMDL) criterion

VS

Standard BIC criterion

