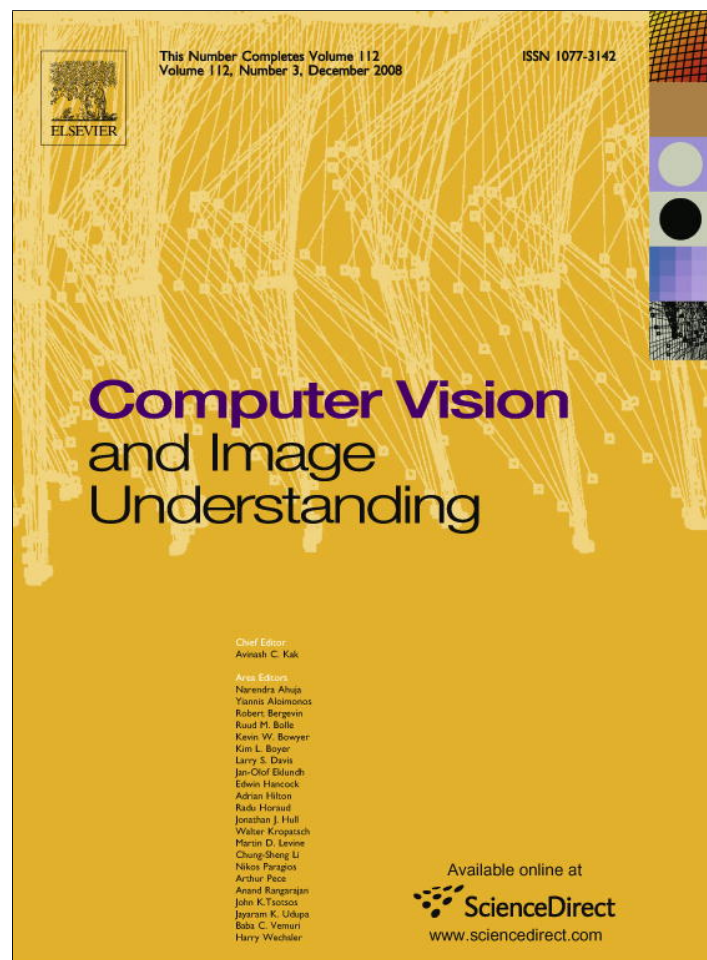


Provided for non-commercial research and education use.
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Contents lists available at ScienceDirect

Computer Vision and Image Understanding

journal homepage: www.elsevier.com/locate/cviu

Optimising dynamic graphical models for video content analysis

Tao Xiang*, Shaogang Gong

Department of Computer Science, Queen Mary, University of London, London E1 4NS, UK

ARTICLE INFO

Article history:

Received 18 July 2006

Accepted 25 May 2008

Available online 8 June 2008

Keywords:

Video content analysis

Structure scoring

Graphical models

Hidden Markov models

Surveillance video segmentation

Group activity modelling

ABSTRACT

A key problem in video content analysis using dynamic graphical models is to learn a suitable model structure given observed visual data. We propose a completed likelihood AIC (CL-AIC) scoring function for solving the problem. CL-AIC differs from existing scoring functions in that it aims to optimise *explicitly* both the explanation and prediction capabilities of a model simultaneously. CL-AIC is derived as a general scoring function suitable for both static and dynamic graphical models with hidden variables. In particular, we formulate CL-AIC for determining the number of hidden states for a hidden Markov model (HMM) and the topology of a dynamically multi-linked HMM (DML-HMM). The effectiveness of CL-AIC on learning the optimal structure of a dynamic graphical model especially given sparse and noisy visual data is shown through comparative experiments against existing scoring functions including Bayesian information criterion (BIC), Akaike's information criterion (AIC), integrated completed likelihood (ICL), and variational Bayesian (VB). We demonstrate that CL-AIC is superior to the other scoring functions in building dynamic graphical models for solving two challenging problems in video content analysis: (1) content based surveillance video segmentation and (2) discovering causal/temporal relationships among visual events for group activity modelling.

© 2008 Elsevier Inc. All rights reserved.

1. Introduction

Dynamic graphical models or dynamic Bayesian networks (DBNs), especially hidden Markov models (HMMs) and their variants, have become increasingly popular for modelling and analysing dynamic video content [12,21,8,11,26,19,15,34]. By using a DBN for video content analysis, we assume that dynamic visual content is generated sequentially by some hidden states of the dynamic scene which evolve over time. These hidden states often have physical meanings. For instance, they could correspond to certain stages/phases of an activity [11,26,21], the occurrence of different classes/phases of visual events [19], or different types of transition segments between video shots [8]. The hidden states, as suggested by the name, cannot be observed directly. They can only be inferred from the observed visual data given a learned DBN. Learning a DBN involves estimating both its structure and parameters from data. The structure of a DBN refers primarily to (1) the number of hidden states of each hidden variables of a model and (2) the conditional independence structure of a model, i.e., factorisation of the state space for determining the topology of a graph. There have been extensive studies in the machine learning community on efficient parameter learning when the structure of the model is known a priori (i.e., assumed) [18]. However, much less effort has been

made to tackle the more challenging problem of learning the optimal structure of an unknown DBN [5,17,13,36]. Most previous DBNs-based video content modelling approaches avoid the structure learning problem by setting the structure manually [21,26,8,15]. However, it has been shown that a learned structure can be advantageous over manually set ones given sparse and noisy visual data [19]. In this paper, we address the problem of how to accurately and robustly learn the optimal structure of a DBN for video content analysis in a realistic situation where only sparse and noisy visual data are available.

Most previous structure learning techniques have adopted a search-and-score paradigm [17].¹ These techniques first define a scoring function/model selection criterion consisting of a maximum likelihood term and a penalty term to penalise complex models. The model structure space is then searched to find the optimal model structure with the highest score. The most commonly used scoring functions include Bayesian information criterion (BIC) [29], minimum description length (MDL) [28], BDe [20], Akaike's information criterion (AIC) [1], integrated completed likelihood (ICL) [6], and variational Bayesian (VB) [5,4]. The selected models are 'optimal' in the sense that they can either best explain the existing data (BIC, MDL), or best predict unseen data (AIC). It has been demonstrated both theoretically and experimentally in the case of static models that explanation oriented scoring functions suffer from model

* Corresponding author. Fax: +44 20 8980 6533.

E-mail addresses: txiang@dcs.qmul.ac.uk (T. Xiang), sgg@dcs.qmul.ac.uk (S. Gong).¹ Alternatives include the Bayesian approach to model selection [5] and context-specific independence [9].

under-fitting while prediction oriented ones suffer from model over-fitting [23,30,6,35].

To address the problems associated with existing scoring functions, we argue that a better scoring function should select a model structure that is capable of both explaining the observed data and predicting unseen data optimally at the same time. To this end, we derive completed likelihood AIC (CL-AIC) for learning the structure of a DBN. CL-AIC was first introduced in our previous work [35] for Gaussian mixture models (GMMs) which can be represented as a static graphical model (see Fig. 1(a)). In this paper, CL-AIC is derived as a general scoring function suitable for both static and dynamic graphical models, with GMMs and DBNs as special cases. In particular, CL-AIC is formulated for determining the number of hidden states of a HMM and for learning the topology of a dynamically multi-linked HMM (DML-HMM) (see Fig. 1(b) and (c)).

The effectiveness of CL-AIC on DBNs structure learning is demonstrated through comparative experiments against BIC, AIC, ICL, and VB. Experiments on synthetic data were carried out to examine and quantify the effect of sample size on the performance of different score functions. The results, for the first time, reveal a key difference in structure learning of static and dynamic graphical models in terms of the definition of data sparseness. We further considered two video content analysis problems using real data: (1) content based surveillance video segmentation and (2) discovering causal/temporal relationships among visual events for group activity modelling. Our experimental results demonstrate that CL-AIC is superior to alternative scoring functions in building dynamic graphical models for video content analysis especially given sparse and noisy data.

The rest of the paper is structured as follows: in Section 2, we derive CL-AIC as a general scoring function for graphical models with hidden variables. We also formulate CL-AIC for two special cases of DBN, namely a HMM and a DML-HMM, and present synthetic experiments to compare CL-AIC to existing scoring functions including BIC, AIC, ICL, and VB. In Section 3, we address the problem of learning the optimal number of video segments for surveillance video segmentation. Comparative experiments are conducted using over 10h of challenging outdoor surveillance video footages. We then compare CL-AIC with other competing scoring functions in learning the topology of a DML-HMM for group activity modelling in Section 4. The paper concludes in Section 5.

2. Completed likelihood AIC for graphical models with hidden variables

We derive CL-AIC for graphical models with hidden variables with GMMs and DBNs as special cases. Let us first consider the nature of computation in estimating and using a graphical model. Consider an observed data set \mathcal{Y} modelled by a graphical model

\mathcal{M}_K with hidden variables. \mathcal{M}_K can be used to perform three tasks: (1) estimating the unknown distribution that most likely generates \mathcal{Y} , (2) inferring the values of hidden variable in \mathcal{M}_K from \mathcal{Y} , and (3) predicting unseen data. Computing (1) and (2) emphasises data explanation while solving (3) concerns with data prediction and synthesising. In this context, scoring functions based on approximating the Bayesian Model Selection criterion [27] such as BIC and VB choose a model that maximises $p(\mathcal{Y}|\mathcal{M}_K)$, the probability of observing \mathcal{Y} given \mathcal{M}_K or the marginal likelihood of the model. They thus enforce mainly task (1). AIC, on the other hand, chooses the model that best predicts unseen data, therefore optimising (3). CL-AIC utilises CL in order to make explicit the task (2), while following a derivation procedure similar to that of AIC.

Completed likelihood was originally derived for as a scoring function for mixture models [7,6]. The model that best explains both the observed data and hidden variables inferred from the observed data is selected by CL as the optimal model. CL is thus a strongly explanation-oriented scoring function. Here, we first extend the definition of CL to a general case of graphical models with hidden variables. The complete data, denoted as $\bar{\mathcal{Y}}$, for such a model is a combination of the observed data (\mathcal{Y}) and the values of the hidden variables (\mathcal{Z}): ($\bar{\mathcal{Y}} = \{\mathcal{Y}, \mathcal{Z}\}$, where \mathcal{Z} is unknown, and must be inferred from \mathcal{Y}). The completed log-likelihood of $\bar{\mathcal{Y}}$ is

$$CL(K) = \log p(\mathcal{Y}|\mathcal{M}_K, \theta_{\mathcal{M}_K}) + \log p(\mathcal{Z}|\mathcal{Y}, \mathcal{M}_K, \theta_{\mathcal{M}_K}),$$

where $\theta_{\mathcal{M}_K}$ are the true model parameters and K is the index of the candidate models. In practice, $\theta_{\mathcal{M}_K}$ are replaced using the Maximum Likelihood (ML) estimate $\hat{\theta}_{\mathcal{M}_K}$ and the unknown values of the hidden variables \mathcal{Z} is replaced by $\hat{\mathcal{Z}}$, the values inferred from the observed data \mathcal{Y} . The completed log-likelihood is thus rewritten as

$$CL(K) = \log p(\mathcal{Y}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}) + \log p(\hat{\mathcal{Z}}|\mathcal{Y}, \mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}). \quad (1)$$

CL-AIC aims to choose a model that best explains the observed data and has the minimal divergence to the true model, therefore best predicting the unseen data as well. The divergence between a candidate model and the true model is measured using the Kullback–Leibler information [24]. Given a completed data set $\bar{\mathcal{Y}}$, we assume that $\bar{\mathcal{Y}}$ is generated by the unknown true model \mathcal{M}_0 with model parameter $\theta_{\mathcal{M}_0}$. For any given model \mathcal{M}_K and the maximum likelihood estimate $\hat{\theta}_{\mathcal{M}_K}$, the Kullback–Leibler divergence between the two models is computed as

$$d(\mathcal{M}_0, \mathcal{M}_K) = E \left[\log \left(\frac{p(\bar{\mathcal{Y}}|\mathcal{M}_0, \theta_{\mathcal{M}_0})}{p(\bar{\mathcal{Y}}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K})} \right) \right]. \quad (2)$$

Ranking the candidate models according to $d(\mathcal{M}_0, \mathcal{M}_K)$ is equivalent to ranking them according to

$$\delta(\mathcal{M}_0, \mathcal{M}_K) = E \left[-2 \log p(\bar{\mathcal{Y}}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}) \right].$$

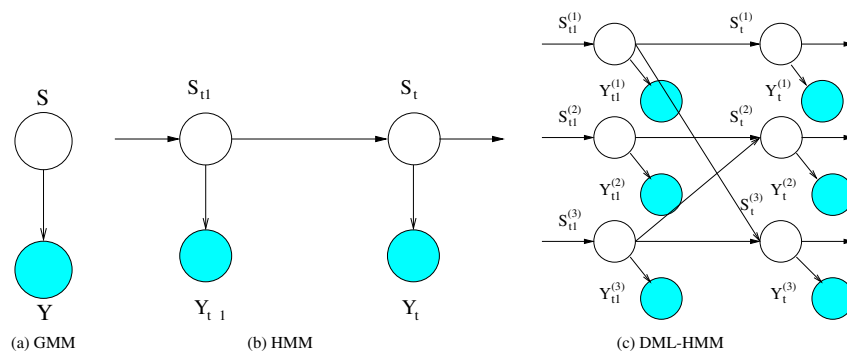


Fig. 1. Three different types of graphical models with hidden nodes, among which HMM and DML-HMM are DBNs. Observation nodes are shown as shaded circles and hidden nodes as clear circles.

$\delta(\mathcal{M}_0, \mathcal{M}_K)$ cannot be computed directly since the unknown true model is required. However, it was noted by Akaike [1] that $-2 \log p(\tilde{\mathcal{Y}}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K})$ can serve as a biased approximation of $\delta(\mathcal{M}_0, \mathcal{M}_K)$, and the bias adjustment

$$E \left[\delta(\mathcal{M}_0, \mathcal{M}_K) + 2 \log p(\tilde{\mathcal{Y}}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}) \right]$$

converges to $2C_K$ when the number of data sample approaches infinity. Our CL-AIC is thus derived as

$$CL - AIC(K) = -\log p(\tilde{\mathcal{Y}}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}) + C_K, \quad (3)$$

where C_K is the dimensionality of the parameter space of model \mathcal{M}_K . The first term on the right-hand side of (3) is the completed likelihood given by Eq. (1). We thus have

$$CL - AIC(K) = -\log p(\mathcal{Y}|\mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}) - \log p(\hat{\mathcal{Z}}|\mathcal{Y}, \mathcal{M}_K, \hat{\theta}_{\mathcal{M}_K}) + C_K. \quad (4)$$

The candidate model structure that has the minimal CL-AIC value will be chosen as the optimal structure.

By utilising an explanation-oriented scoring function (CL) and following a derivation procedure similar to that of a prediction-oriented scoring function (AIC), CL-AIC is able to select the model structure that maximises both the probability of the observed data and the posterior probability of the inferred hidden/incomplete data simultaneously whilst having a minimum number of free parameters (i.e., as simple as possible). This formulation results in a number of important differences compared to existing techniques: (1) Unlike previous scoring functions, CL-AIC attempts to optimise *explicitly* the explanation and prediction capabilities of a model. This makes CL-AIC theoretically attractive. The effectiveness of CL-AIC in practice is demonstrated through experiments in the following sections. (2) It has been shown that CL can be combined with BIC which leads to an ICL scoring function [6]. However, the experiments reported in [6] indicated that in the case of mixture models, ICL performs poorly when data belonging to different mixture components are severely overlapped. We suggest this is caused by the factor that ICL is a combination of two explanation-oriented scoring functions without considering the prediction capability of the model. Since CL-AIC integrates an explanation-oriented scoring function with a prediction-oriented one, it is theoretically better justified than ICL. Our experiments in the following reinforce this observation.

As shown above, CL-AIC is derived based on a very intuitive principle and can be formulated for a general category of models, namely graphical models with hidden variables. Example of these models include both static statistical models such as GMMs, naive Bayesian classifiers, principal component analysis (PCA), and Independent Components Analysis (ICA), and dynamic ones such as HMMs and their variants, and Kalman Filters. Since we are interested in Dynamic Bayesian Networks (DBNs) for video content analysis in this paper, we formulate CL-AIC for two cases of DBNs in the following.

2.1. Determining the number of hidden states of a HMM

Let us now consider a specific problem of learning the structure of a hidden Markov model (HMM). A HMM has a fixed topology with one hidden variable and one observation variable at each time instance t (see Fig. 1(b)). The hidden variable is discrete in most applications. The structure learning problem for a HMM thus refers to how to determine the number of hidden states that the hidden variable can assume. Assuming that at each time instance t , the discrete hidden variable S_t can have K different values (states), the complete data for the model is $\tilde{\mathcal{Y}} = \{\mathcal{Y}, \mathcal{Z}\}$ where \mathcal{Z} is the true hidden variable values over different time instances (i.e. the true hidden state sequence). The completed log-likelihood of $\tilde{\mathcal{Y}}$ is computed as

$$CL(K) = \log \left(\sum_S p(\mathcal{Y}|S, \hat{\theta}_K) p(S|\hat{\theta}_K) \right) + \log p(S = \hat{\mathcal{Z}}|\mathcal{Y}, \hat{\theta}_K), \quad (5)$$

where $S = \{S_1, \dots, S_T\}$ represents all the possible hidden state sequences, T is the length of the sequence, $\hat{\theta}_K$ are the ML (maximum likelihood) estimate of the model parameters of a HMM with K hidden states, $\hat{\mathcal{Z}}$ is the most probable state sequence (i.e. the hidden state sequence among S that best explains the observation sequence) given $\hat{\theta}_K$ and \mathcal{Y} . $\hat{\theta}_K$ can be computed using the forward-backward (Baum-Welch) algorithm [3] based on dynamic programming and $\hat{\mathcal{Z}}$ can be obtained using the Viterbi algorithm [16]. We thus have

$$CL - AIC(K) = -\log \left(\sum_S p(\mathcal{Y}|S, \hat{\theta}_K) p(S|\hat{\theta}_K) \right) - \log p(S = \hat{\mathcal{Z}}|\mathcal{Y}, \hat{\theta}_K) + C_K. \quad (6)$$

A detailed description on how to use the forward-backward procedure and the Viterbi algorithm to efficiently compute the first and second terms on the right hand side of Eq. (6) can be found at Rabiner's excellent HMM tutorial [25]. The optimal number of hidden states \hat{K} is then determined as

$$\hat{K} = \arg \min_K CL - AIC(K). \quad (7)$$

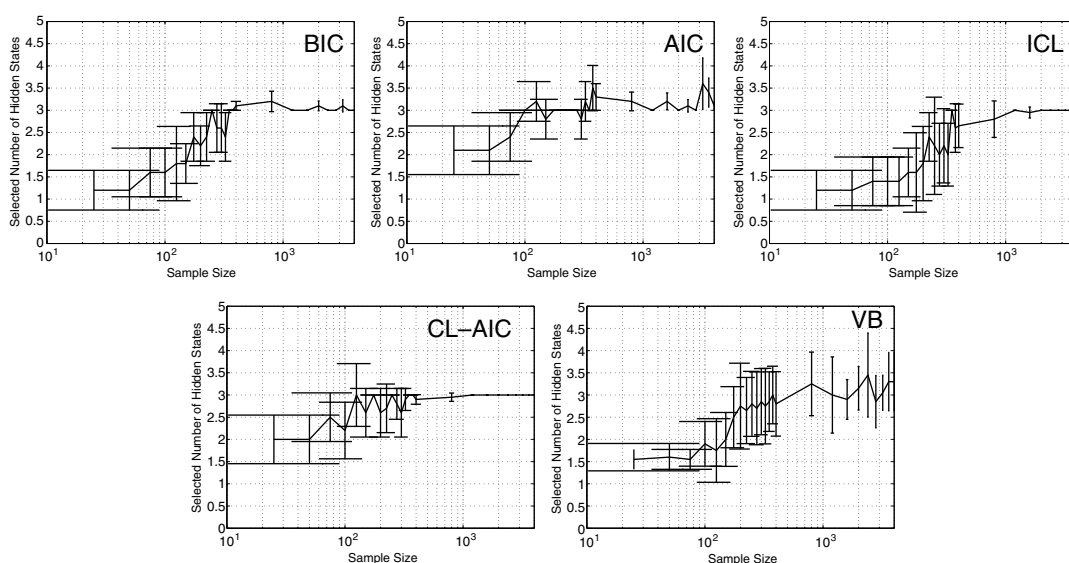
Synthetic experiments were conducted to compare the effectiveness of CL-AIC with that of BIC, AIC, ICL, and VB on determining the number of hidden states of a HMM given data of different sample sizes. One-dimensional data were first generated from a 3-state HMM (i.e., the hidden variable at each time instance can assume 3 states) whose parameters are

$$\mathbf{A} = \begin{bmatrix} 1/3 & 1/6 & 1/2 \\ 0 & 1/3 & 2/3 \\ 1/2 & 1/2 & 0 \end{bmatrix}, \quad \boldsymbol{\pi} = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}, \quad \mathbf{B} = \begin{cases} \mu_1 = 1, \sigma_1^2 = 0.5 \\ \mu_2 = 3, \sigma_2^2 = 0.5 \\ \mu_3 = 5, \sigma_3^2 = 0.5 \end{cases}, \quad (8)$$

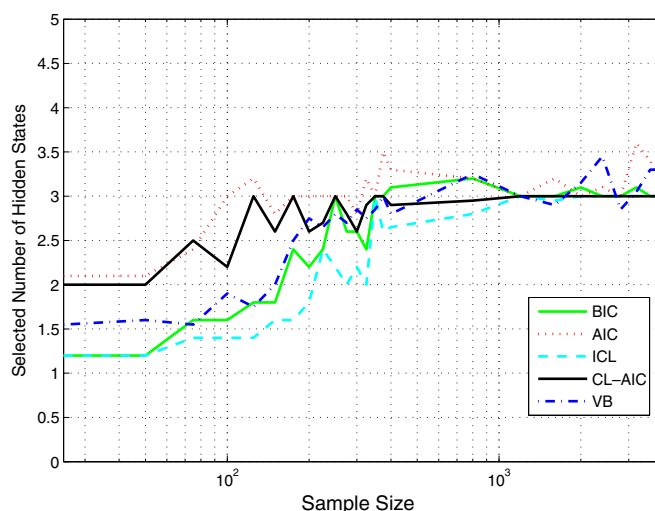
where \mathbf{A} is the transition matrix, $\boldsymbol{\pi}$ is the initial state probability and \mathbf{B} contains the parameters of the emission density (Gaussians with the indicated means and variances). The total number of free parameters is 14 for this HMM. The data were then perturbed by uniformly distributed random noise with a range of $[-0.50, 0.50]$. HMMs with the number of hidden states K varying from 1 to 10 were evaluated. The five different scoring functions were tested on the data sets with the sample size T varying from 25 to 4000. To avoid being trapped at local minima, the EM algorithm² used for estimating model parameters was randomly initialised for 20 times and the solution that yielded the largest likelihood after 30 iterations were chosen. The results are shown in Fig. 2 using the mean and ± 1 standard deviation of the selected number of hidden states over 50 trials, with each trial having a different random number seeds.

Fig. 2(b) shows the mean of the number of states estimated by different scoring functions over 50 trials in a single plot. It can be seen that when the sample sizes were small, all five scoring functions tended to favour under-fitted models (i.e., underestimate the number of states), with AIC and CL-AIC clearly outperforming BIC, ICL and VB. As the sample sizes increased, the number of hidden states determined using all scoring functions converged to the true number 3. Given densely sampled data sets ($T > 400$), our results show that both AIC, BIC, and VB tended to slightly over-fit while ICL and CL-AIC yielded more accurate estimation of the number of hidden states. Fig. 2(a) shows the variations of the structure learning results across different trials which indicate the sensitiv-

² A modified EM algorithm was used for computing the VB scoring function [4].



(a)



(b)

Fig. 2. Synthetic data experiment results for determining the number of hidden states of a HMM using different scoring functions. (a) Selected number of hidden states (mean and ± 1 standard deviation over 50 trials). (b) Mean of the selected number of hidden states (the true number of hidden states is 3).

ity of different scoring function to the distribution of data and more importantly the presence of noise. It can be seen from Fig. 2(a) that both AIC and VB exhibited large variations in the estimated number of states no matter what the sample size was, while other scoring functions had smaller variations given larger sample sizes. Fig. 3 shows two examples of determining the optimal number of hidden states using different scoring functions. The value of different scoring functions is plotted against the number of hidden states in the candidate models. It can be seen that the scoring function value plots of AIC and VB have more local minimal and are less discriminative for different candidate model structures than the other three. This partially explains why the structure learning results for AIC and VB exhibited larger variations across different trials. Overall, the results show that CL-AIC has the best performance among the five scoring functions that we compared.

It can be seen from Figs. 2 and 3 that different scoring function perform differently as the sample size increases. In particular, Fig. 2(b) shows that when the data set is sparse, all five scoring functions tend to underfit the model. In the experiment presented in Figs. 2 and 3, the distribution of the observed data follows three

Gaussians (see Eq. (8)) and each hidden state will correspond to one Gaussian when both the model structure and model parameters are learned accurately. When the sample size is small with the presence of noise, some of the three Gaussians may not be well supported by the data. This explains the underfitting tendency of the five scoring functions. Both Figs. 2(b) and 3(a) also show that given a sparsely sampled data set, CL-AIC and AIC outperform the other three. It is easy to understand why CL-AIC is less underfitting than ICL in this case. A comparison of the two shows that they only differ in their penalty terms in formulation. In particular, CL-AIC has a weaker penalty term, therefore always favouring more complicated model compared to ICL. For exactly the same reason, AIC is less underfitting than BIC in Figs. 2(b) and 3(a) given sparse data. Both Figs. 2(b) and 3(a) also suggest that the hidden data likelihood terms in CL-AIC and ICL played an insignificant role when the size of the data set is small. This is reflected by the fact that CL-AIC and AIC behaved very similarly given sparse data, so did ICL and BIC.

When the data set is densely populated, the three Gaussians are all well supported by the data samples. Therefore, all scoring functions have a much easier task to estimate the number of hidden

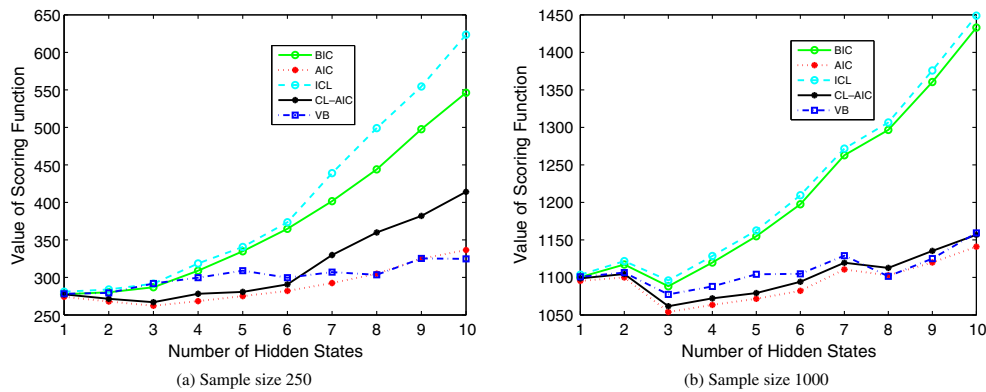


Fig. 3. Examples of choosing optimal number of hidden states using different scoring functions. The value of five different scoring functions are plotted against the number of hidden states. The global minimal on the plots correspond to the selected numbers of hidden states. In (a), both AIC and CL-AIC chose the correct hidden state number 3 while the other three select 2 given a data set of 250 samples. In (b) all five scoring functions selected the correct number given a data set of 1000 samples.

states. However, with the randomly distributed noise, the distribution of the data contained in each original Gaussian can become non-Gaussian. We know that a non-Gaussian distribution can be approximated using multiple Gaussians, which will then result in more than three hidden states in the estimated model structure. This explains why the models selected by VB, BIC and AIC are slightly overfitted. Nevertheless, Fig. 2 shows that CL-AIC and ICL do not suffer from the same problem. This is because for both of them, complete likelihood is part of the scoring function, which means that an estimated model must explain explicitly the hidden state sequence as well. The temporal order of the hidden states that correspond to the noise-caused extra Gaussians cannot be well explained by the data. For instance, if two Gaussians are now used to approximate an original single Gaussian, the temporal order of these two Gaussians will be random. This results in low likelihood of the hidden data. In other words, a model that is more complicated than the true one would be penalised by both the penalty term and the term corresponding to the likelihood of the hidden data in CL-AIC and ICL. Therefore, both ICL and CL-AIC give more accurate estimates when sample size is large given noisy data compared to the other three scoring functions.

Both Figs. 2 and 3 show that VB is sensitive to initialisation even with densely sampled data (see the error bar in Fig. 2(a) and the local minima on the plots corresponding to VB in Fig. 3). Note that the key difference between computing VB and the other four scoring functions is that the former is based a modified EM algorithm and the latter use standard EM algorithm [4]. The modified EM algorithm involves computing a distribution over the model parameters rather than a point estimate of them in the M-step. This makes VB even more sensitive to initialisation than the other four. Consequently, both the plots in Fig. 3(a) and (b) that correspond to VB are less smooth with more local minima than the others.

2.2. Determining the topology of a DML-HMM

We now consider the problem of determining the unknown topology of a dynamically multi-linked HMM (DML-HMM) [19] from data using CL-AIC as the scoring function. DML-HMMs were first proposed in [19] to overcome a limitation of HMMs in modelling complex dynamic processes. In particular, a HMM requires a large number of parameters to describe if it is to model multiple temporal processes simultaneously. This implies that a single hidden state variable and a single observation variable are to represent implicitly multiple sources of variations at any given time instance. Unless the training data set is very large and relatively ‘clean’, poor model learning is expected. To address this problem, various topological extensions to the standard HMMs can be considered to fac-

torise the observation and/or state space by introducing multiple hidden state variables and multiple observation variables for modelling different temporal processes explicitly and simultaneously. Instead of being fully connected as in the case of a coupled HMM (CHMM) [12] or being without any inter-links between different temporal processes as in the case of a parallel HMM (PaHMM) [32], a DML-HMM aims to *only* connect a subset of relevant hidden state variables across multiple temporal processes. This is in order to learn the optimal factorisation of the state space and more importantly to be less sensitive to data noise. Instead of being fixed as in the cases of CHMMs and PaHMMs, the structure of a DML-HMM is learned from data. Given a data set, it is assumed that at each time instance the temporal process responsible for each data sample is known and the number of hidden states for each hidden variable is also known. The unknown structure to be learned is thus the topology of the graph, i.e., the links among different hidden nodes within two consecutive time instances.

To learn the optimal topology of a DML-HMM, CL-AIC is computed for each candidate topology using Eq. (6) where the subscript K becomes the index of different topologies. The first and second terms on the right-hand side of Eq. (6) can be computed using an extended forward-backward algorithm for a DML-HMM. A detailed description of the algorithm can be found in the appendix of this paper. Note that the total number of candidate topologies K_{max} is exponential in the number of temporal processes N_t . Each candidate topology can be represented using a $N_t \times N_t$ inter-connection matrix whose elements have value 1 if there is a directed link between the two corresponding hidden nodes within two consecutive time instances and 0 otherwise.

Synthetic experiments were carried out to examine the effectiveness of CL-AIC in comparison with that of BIC, AIC, ICL on determining the topology of a DML-HMM given data of different sample sizes.³ Data sets were randomly generated using a DML-HMM with three temporal processes whose topology is shown in Fig. 1(c). The model inter-connection matrix is thus $\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$. Both the observation and hidden variables are discrete with three possible values. The total number of free parameters of the model is therefore 66. The model parameters were set randomly. One-fourth of the observational data were replaced by random numbers to synthesise noise contained in the observation. DML-HMMs with $K_{max} = 64$ different topologies were evaluated by four different scoring functions using

³ Variational Bayesian (VB) was not evaluated here because implementing VB for a dynamic graphical model with multiple temporal processes such as a DML-HMM is not trivial and beyond the scope of this paper.

data sets with sample size T varying from 25 to 4000. The performance of different scoring functions was measured by looking at both the number of links connecting hidden nodes within two consecutive time instances (the true number is 5) and the hamming distance between the estimated inter-connection matrices and the true one (the distance is zero if the structure is estimated correctly). The former measures model complexity while the latter measures the accuracy of the learned structure. The experimental results were obtained over 50 trials, and are shown in Fig. 4.

It can be seen from Fig. 4 that given sparse and noisy data the optimal models selected using all four different scoring functions tended to underfit with ICL and CL-AIC outperforming the other two. As the sample sizes increased, the optimal number of links among hidden nodes selected by CL-AIC and ICL converged towards the true number 5, while those selected by BIC and AIC converged to 4, (i.e., underfitting). In the meantime, the hamming distance obtained using different scoring functions decreased, with that obtained using CL-AIC being the smallest. Our experiment shows that all four scoring functions are sensitive to initialisation

of the model parameters (see the large standard deviations depicted in Fig. 4(a) and (c)). The main reason for this shortcoming is that the EM algorithm used for learning model parameters and estimating the likelihood terms in different scoring functions is sensitive to initialisation. This is well reported in the literature. Although we have taken measures to avoid this problem by randomly initialising the EM algorithm and choosing a solution that yields the highest likelihood, the problem is not totally eliminated. This suggests that each scoring function has a tendency to form a complicated surface with a large number of local minima with respect to the candidate model topologies. To solve this problem, more sophisticated method needs be exploited/developed.

Apart from this common limitation, the results also show that BIC and AIC tend to underfit the model whilst CL-AIC and ICL do not suffer from that problem given sufficient data. There is an intuitive explanation for this. Essentially for both BIC and AIC, the likelihood of the observed data would not benefit directly from having interlinks between different time processes. In the meantime, the penalty terms in their formulation would increase exponentially

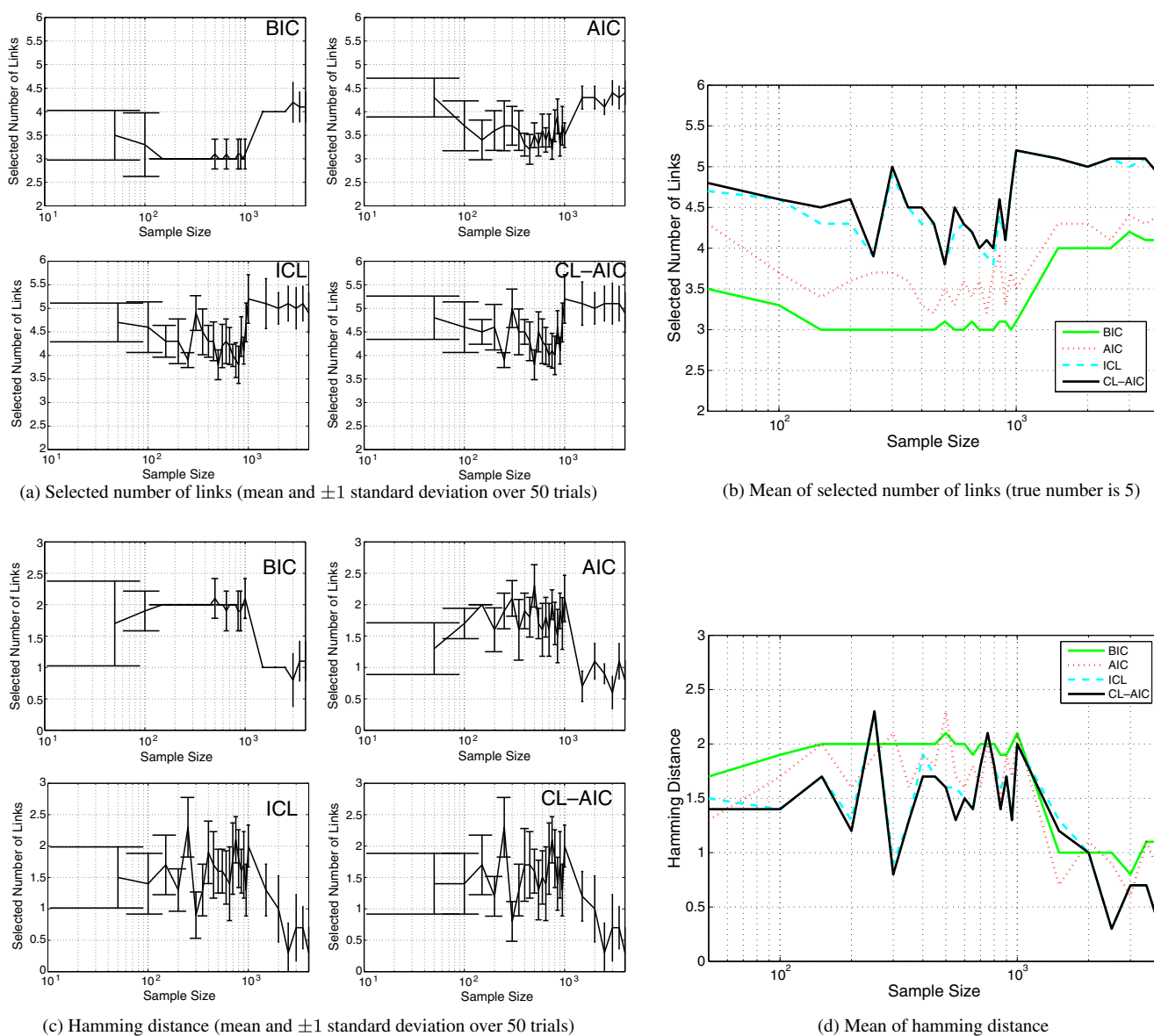


Fig. 4. Synthetic data experiment results for determining the topology of a DML-HMM using different scoring functions. In (a) and (b), the true number of links is 5. In (c) and (d), smaller Hamming distances indicate more accurate structure learning.

with the increase of the number of interlinks. It is therefore natural to observe that underfitted models are preferred by BIC and AIC even given sufficient data. As for CL-AIC and ICL, they do not have the same problem because the likelihood of the inferred hidden state sequence is included in their formulation. Having multiple interlinks, especially those correct ones, helps greatly in explaining away the temporal occurrence of different states and therefore leads to higher likelihood of the hidden data.

3. Surveillance video segmentation

HMMs have been widely used for automatic segmentation of sequential/time-series data such as speech [14], DNA sequences [10] and video [8,11,34]. Here we propose to use HMM for content based surveillance video segmentation, i.e. to segment a continuous surveillance video based on activities captured in the video. Note that since there is only one video shot in a continuous surveillance video, the conventional shot-change detection based segmentation approach [2] cannot be adopted. We thus employed a video content trajectory based method. Specifically, a L -dimensional feature vector is first extracted from each image frame. The video content is thus represented a video trajectory in this L -dimensional feature space. This feature vector is then represented as the observational input of a HMM at each time instance. The conditional probability distributions (CPDs) of each observation variable are Gaussian for each of the K states of its parent hidden variable. The video content is then monitored using the discrete hidden variables in the model. The changes of video content can thus be detected as the changes of hidden states which correspond to breakpoints on a video trajectory (N detected change points/breakpoints result in $N + 1$ video segments for a continuous video). Using a left-to-right HMM model [25], the number of hidden states would correspond to the number of video segments which can be automatically determined using CL-AIC.

3.1. Activity-based video content representation

For clarity and concreteness, we shall illustrate our approach using surveillance videos monitoring aircraft docking operations at an airport. The total number of activities taking place in the aircraft docking area was decided as 15 by airport ground control officers. A trained ground control officer in the control room is required to monitor these 15 activities via both video inputs from multiple cameras and audio communications with various staff members on the ground. However, this does not imply that all of the 15 activities to be monitored by ground control officers are identifiable when CCTV video is the only input. It turns out that, when watching video for identifying these 15 activities on the list, most of them are completely or partially blocked by the aircraft

and the air bridge (e.g., crew bus picking up crew, engineers doing service check) and only six activities on the list are visually detectable. This is because the data were captured using a single camera from the frontal view. Obviously if more cameras were deployed to cover different angles, the number of detectable activities will increase. But it is still unlikely that all 15 activities will be visually detectable due to the limit of the cameras one can possibly install in reality and also the long distance from the camera to the objects of interest resulting in low image resolution.

Fig. 5 shows an example activity based structure of a complete aircraft docking operation. One activity can be followed by either another activity immediately or a period of inactivity when no meaningful visual changes can be detected. The durations of activities and inactivity gaps vary hugely during a single aircraft docking operation and across different operations. There are also variations in the temporal order of activities and inactivity gaps. For example, the temporal order of `frontalCargoService` and `frontalCateringService` appears to be arbitrary. It is noted that some activities such as `airCRAFTArrival` involve the movement of a single object while other activities such as `frontalCargoService` and `frontalCateringService` consist of the movement of multiple objects which may leave and re-appear in the scene. For the latter case, there often exist a number of short inactivity break-ups within an activity which are different from a long inactivity gap between two activities only by their durations. All these characteristics make analysis of the video content very challenging.

We consider an activity based video content representation [33]. Visual events are detected and classified automatically in the scene. The semantics of video content are considered to be best encoded in the occurrence of such events and the temporal correlations among them. First, an adaptive Gaussian mixture background model [31] is adopted to detect foreground pixels which are modelled using pixel change history (PCH) [33]. Second, the foreground pixels in vicinity are grouped into a blob using the connected component method. Each blob with its average PCH value greater than a threshold is then defined as an event. A detected event is represented as a seven-dimensional feature vector

$$\mathbf{v} = [\bar{x}, \bar{y}, w, h, R_f, M_{px}, M_{py}], \quad (9)$$

where (\bar{x}, \bar{y}) is the centroid of the blob, (w, h) is the blob dimension, R_f is the filling ratio of foreground pixels within the bounding box associated with the blob, and (M_{px}, M_{py}) are a pair of first-order moments of the blob represented by PCH. Among these features, (\bar{x}, \bar{y}) are location features, (w, h) and R_f are principally shape features but also contain some indirect motion information, and (M_{px}, M_{py}) are motion features capturing the direction of object motion. Classification is then performed in the 7D event feature space over a collection of training videos using a Gaussian mixture model (GMM). The number

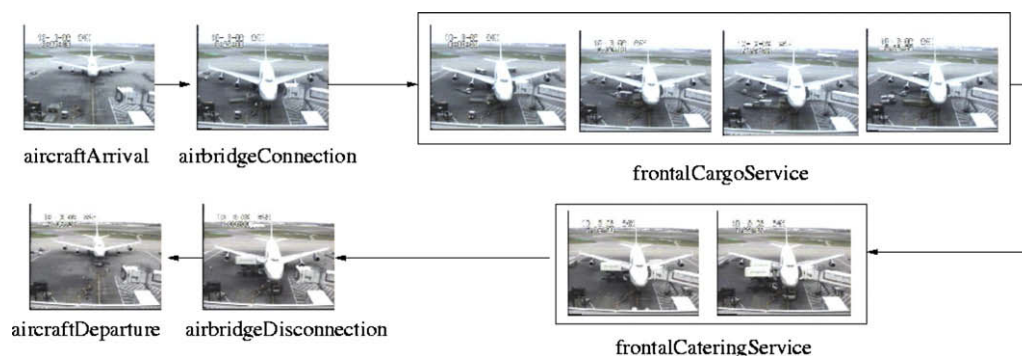


Fig. 5. An example structure of a complete aircraft docking video. Representative frames of different activities occurred in the scene are also shown.

of event classes captured in videos is determined using CL-AIC (see [35] for the formulation of CL-AIC for GMMs). For the aircraft docking scene videos, eight classes of events were automatically detected. It is worth mentioning that different events occur simultaneously and such an event detection mechanism makes mistakes. Mis-detection and wrong labelling can be caused by discontinuous movement and closeness of different objects. An example of event detection and classification is shown in Fig. 6.

Classified events can be considered as ‘snapshots’ of activities captured in a scene. To hide the image feature information and focus on the semantic video content, a scene vector is constructed for each image frame of a video. A scene vector \mathbf{sv}_t for a video frame t is defined as

$$\mathbf{sv}_t = [s_t^1, \dots, s_t^k, \dots, s_t^K], \quad (10)$$

where K is the number of event classes automatically determined using CL-AIC. The value of s_t^k is the number of events of the k th event class detected in the frame t . A scene vector gives a description of ‘what is happening’ in the scene through the class labels of

the detected events. It is thus a concise representation of the video content at the semantic level (see Fig. 7 for an example). However, directly using this scene vector to detect video content changes can cause problems for video segmentation. Specifically, the value of a scene vector \mathbf{sv}_t can become $\mathbf{0}$ (i.e., absent of any event at a given frame) frequently throughout the video sequence (see Fig. 7). This can be caused by either frequent but short inactivity break-ups within activities or long inactivities between activities. Each ‘coming to zero’ is reflected as a dramatic change in the value of \mathbf{sv}_t due to the discrete nature of s_t^k . Those changes that correspond to real changes of video content can thus easily be overwhelmed by changes caused by the inactivity break-ups within activities, which makes temporal segmentation of the video difficult.

To overcome this problem, let us now consider representing the video content using a cumulative scene vector computed at frame t using \mathbf{sv}_t from frame 1 to the frame t . More specifically, the k^{th} element of the cumulative scene vector (denoted as $\tilde{\mathbf{sv}}_t$) is computed as

$$\tilde{s}_t^k = \sum_{i=1}^t s_i^k. \quad (11)$$

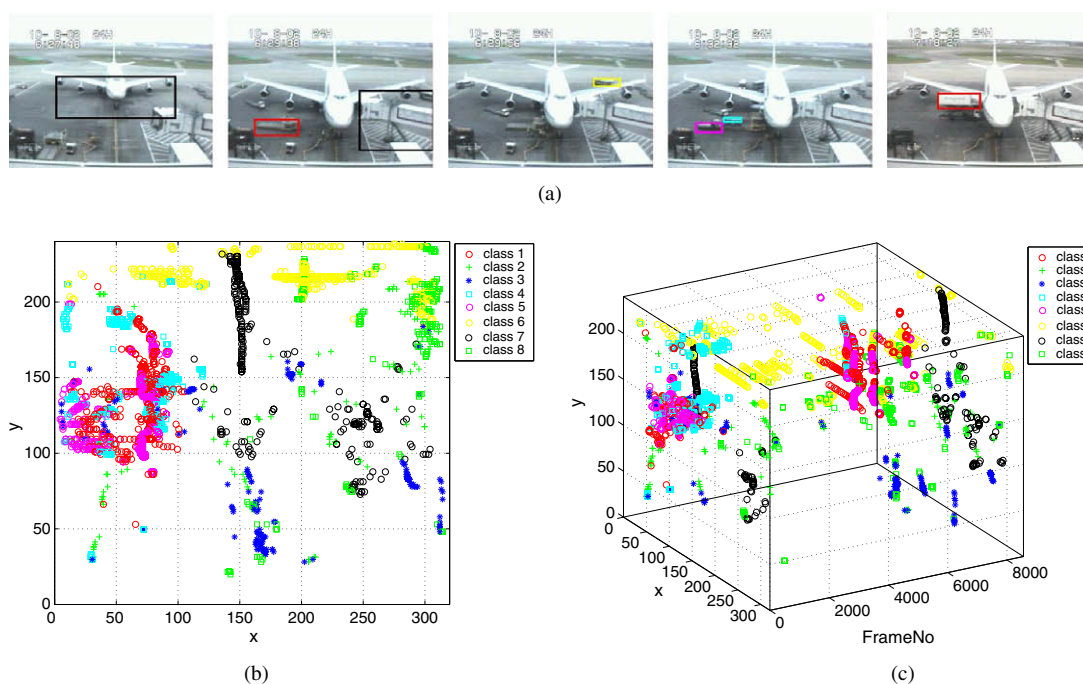


Fig. 6. Event detection and classification in an aircraft docking scene video. Eight classes of events are detected automatically, each of which corresponds to different movement patterns in the image frame. For example, event class 7 corresponds to the movements of the aircraft. Event classes 4, 5, and 6 correspond to the movements of objects involved in activities *frontalCargoService* and *frontalCateringService*. (a) Events are detected and classified into different classes which are highlighted using bounding boxes in different colours. The spatial and temporal distribution of events of different classes are illustrated in (b) and (c), respectively, with centroids of different classes of events depicted using different colours. (For interpretation of color mentioned in this figure the reader is referred to the web version of the article.)

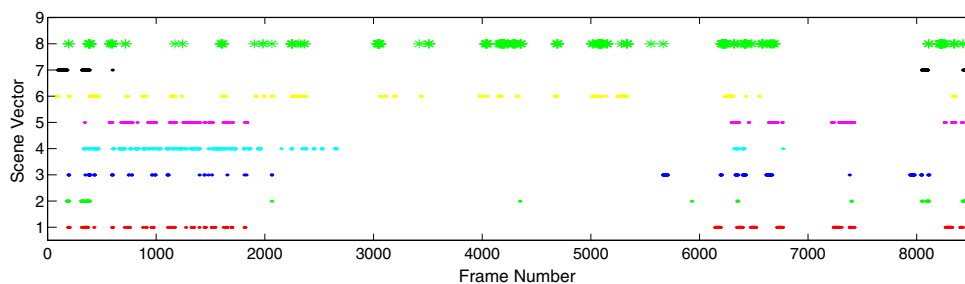


Fig. 7. The example video shown in Fig. 6 is represented by scene vectors evolving over time. We have $K = 8$ for this video. s_t^k is depicted by a dot in colour when $s_t^k > 0$. There are frequent but short inactivity break-ups within activities (see between frame 1 and frame 2000) and a long inactivity gaps between activities (between frame 3000 and frame 6000).

The value of each element of $\tilde{\mathbf{sv}}_t$ becomes continuous and will increase monotonically with time (see Fig. 8). Compared to the scene vector representation \mathbf{sv}_t (see Fig. 7), the short inactivity break-ups at individual frames have little impact on the values of the cumulative scene vector evolved over time. It thus becomes easier to detect breakpoints that correspond to significant changes in video content.

3.2. Determining the number of video segments

Now each image frame of the aircraft docking scene videos is represented as a 8-dimensional feature vector which is then utilised as observational input for a HMM for temporal segmentation. The problem to be solved here is to automatically determine K , the number of hidden states which corresponds to the number of video segments.

Our data set consists of seven aircraft docking scene videos captured using a fixed CCTV analogue camera. After digitalisation, the final video sequences have a frame rate of 2 Hz. Note that it is not uncommon to have such an extremely low frame rate for CCTV surveillance videos, which makes the video segmentation problem even more challenging. Each image frame has a size of 320×240 pixels. The seven video sequences of aircraft docking lasted from 6470 to 17,262 frames per sequence (around 50–140 min of recording), giving in total 72,776 frames (10 h) of video data that cover different times of different days under changing lighting conditions, from early morning, midday to late afternoon. They are referred as video 1 to video 7, respectively. The seven videos were first manually segmented into activities to give the ground truth of the breakpoints for segmentation, resulting in a total of 64 breakpoints and 71 segments. The lengths of these video segments were within the range of 127–3210 frames.

The performance of different score functions are compared by looking at the number of detected breakpoints, the number of true positives, false positives, and false negatives. The results are shown in Table 1 and Fig. 9. Given the true number of breakpoints 64, it can be seen from Table 1 that both BIC, ICL, and VB underestimated

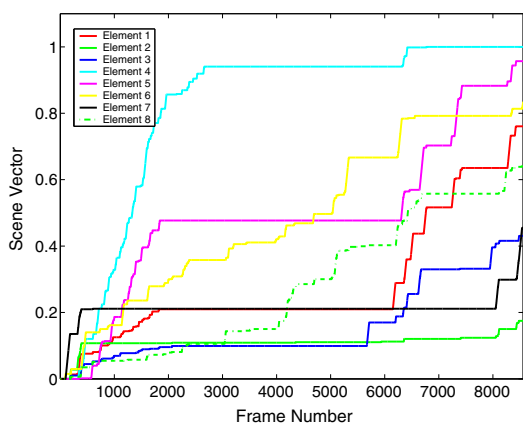


Fig. 8. The eight elements of $\tilde{\mathbf{sv}}_t$ over time for the example video shown in Fig. 6. Each element of $\tilde{\mathbf{sv}}_t$ is normalised to have a value range of [0, 1].

Table 1
Comparing different scoring functions for video segmentation (the true number of breakpoints was 64)

	BIC	AIC	ICL	CL-AIC	VB
# Det. B points	49	73	45	62	52
# True positives	39	52	37	54	40
# False positives	10	21	8	8	12
# False negatives	15	12	17	10	14

the number of segments while AIC overestimated the segment number. In the meantime, the number of segments estimated using CL-AIC was the closest to the true number. On the accuracy of breakpoint detection, Table 1 shows that CL-AIC yielded the highest true positive number and lowest false positive and false negative numbers. In the meantime, both BIC, ICL, and VB gave low false positive number but low true positive number as well. As for AIC, high true positive number was obtained at the price of high false positive number.

4. Discovering causal relationships in group activity modelling

A group activity involves multiple objects co-existing and interacting in a shared common space. Examples of group activities include 'people playing football' and 'shoppers checking out at a supermarket'. Group Activity modelling is concerned with not only modelling actions executed by different objects in isolation, but also the interactions and causal/temporal relationships among these actions. Adopting a DML-HMM based activity modelling approach [19], we consider that a group activity is composed of different classes of dynamically linked visual events representing significant changes in the image over time caused by different objects in the scene. An event is represented by a multi-dimensional feature vector and automatically detected and classified into different event classes (see Section 3.1 for details). The detected events are then taken as the observational input to a DML-HMM so that learning causal and temporal relationships among different classes of events can be achieved by learning the optimal structure of the DML-HMM for modelling the dynamics of the detected events and the interactions among them. More specifically, each temporal process of the DML-HMM is used to model the dynamics of one class of events and those links among different processes capture the causal/temporal relationships of different classes of events. An inter-link between two temporal processes therefore indicates the existence of causality between the two corresponding classes of events.

Experiments were carried out to compare the performance of different scoring functions for learning the structure of a DML-HMM and therefore the causal relationships among different classes of visual events. A simulated 'shopping scenario' was captured on a 20 min video. Some typical scenes can be seen in Fig. 10(a). The scene consists of a shopkeeper sitting behind a table on the right side of the view. Drink cans were laid out on a display table. Shoppers entered from the left and either browsed without paying or took a can and paid for it. The data used for this experiment were sampled at 5 frames per second with total number of 5699 frames of images sized 320×240 pixels.

In the 20-min video, a total of 4634 events were automatically detected and classified into five event classes. By observation, the five classes of events corresponds to five key constituents of the shopping activity. They were labelled as can-Taken, entering/leaving, shopkeeper, browsing and paying respectively (see Fig. 10(a)). It was noted that different classes of events occurred simultaneously. It is also true that our event recognition model made errors. Some of the errors were caused by the occlusion, closeness and visual similarity among different events. Some others were due to the factor that the causal/temporal relationships among events were not considered at the level of event detection. For example, when a shopper stands in front of the shopkeeper, it is impossible to tell whether he/she is going to pay unless one takes into consideration whether any drink can was taken a moment ago. The event classifier is therefore expected to make such errors without taking into account the temporal and causal correlations among different classes of events. Such causal/temporal relationships are modelled using a DML-HMM.



(a) Different activities captured in video 2

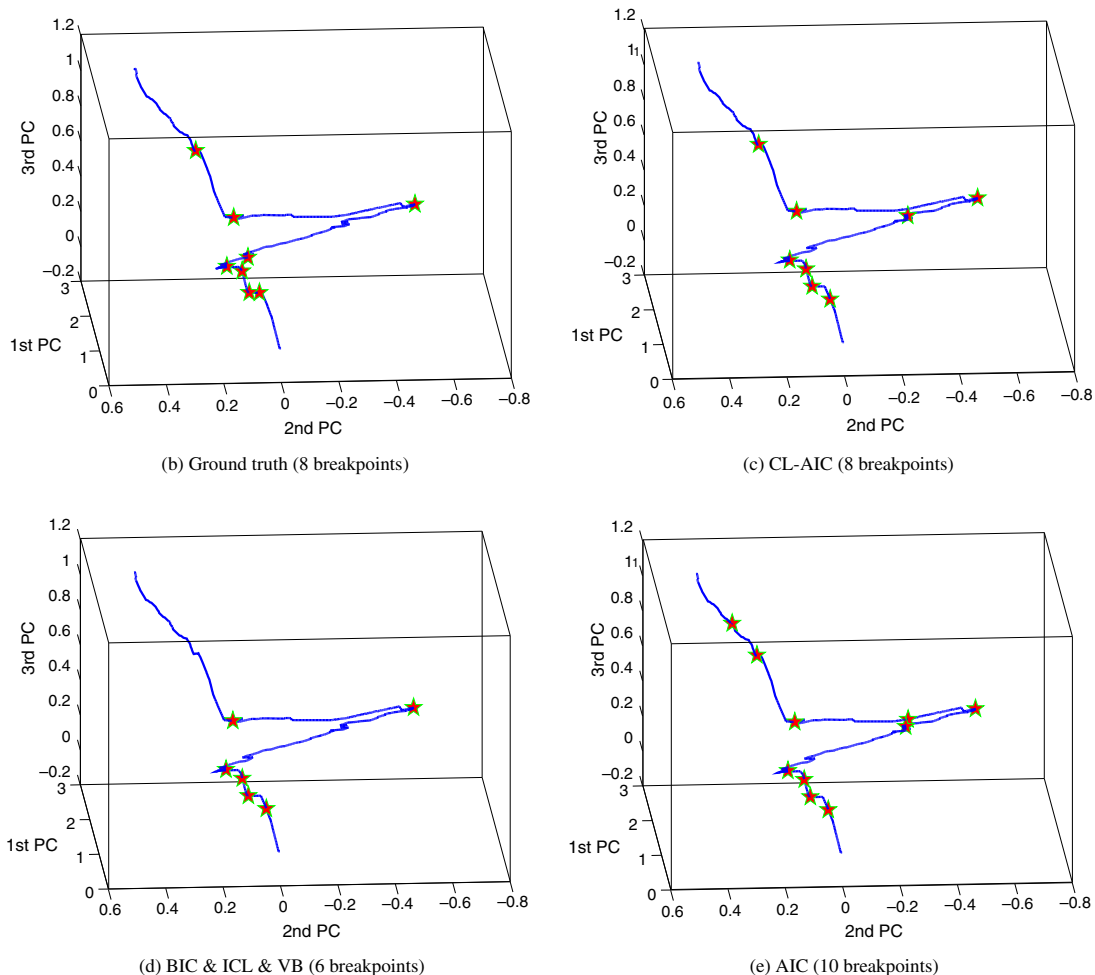


Fig. 9. Determining the number of video segments for an aircraft docking video (video 2 of the 7 videos) using a HMM with different score functions. (a) Representative frames of different activities captured video 2. These activities were (from left to right): aircraftArrival, airbridgeConnection, frontalCargoService, frontalCateringService, airbridgeDisconnection, and aircraftDeparture. (b) Ground truth obtained by manually segmenting the video. (c)–(e) segmentation results using different scoring functions with the detected breakpoints shown on the video trajectory. Note that in (b)–(e) the video trajectories are shown in a 3D PCA space of the original 8D video content feature space just for the illustration purpose (the computation was done in the original 8D feature space).

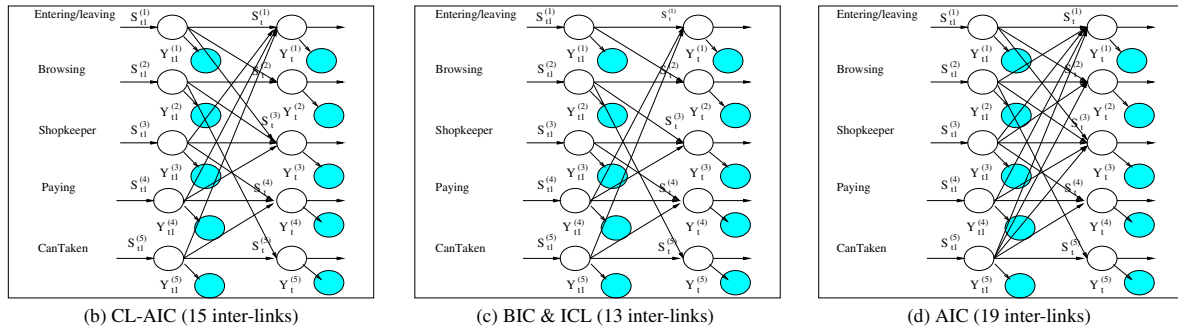
There are five temporal processes in this DML-HMM, each corresponding to one class of events. We also consider two states for each hidden state variable, i.e., a binary variable switching between the status of *True* and *False* which correspond to whether or not an event of a certain class is truly present in each frame. Each observation variable is continuous and given by a 7D feature vector representing a detected event (see Eq. (9)). Their distributions are mixtures of Gaussian with respect to the states of their discrete parent nodes. For model learning, the distributions of the detected events learned using GMMs are used to initialise the distributions of the observation vectors. The priors and transition matrices of states are initialised randomly. The number of candidate topologies for a 5-temporal-process DML-HMM is too large to be searched exhaustively. The Structural EM algorithm [17] was thus adopted to search for the optimal structure more efficiently using different scoring functions.

The discovered causal/temporal relationships among different classes of events are embedded in the learned topologies of the DML-HMMs. For instance, a link pointing from the *canTaken* process towards the *paying* process indicates the causality between these two classes of events. Compared with the expected structure of the shopping activity shown in Fig. 10(e), it can be seen that the causal relationships among different classes of events and the temporal structure of the activity were discovered correctly by CL-AIC (Fig. 10(b)). In comparison, an over-complicated DML-HMM topology was selected using AIC (Fig. 10(d)) while both BIC and ICL underestimated the number of inter-links among different temporal processes (Fig. 10(c)), resulting in over-simplified causal relationships among different classes of events.

Note that in Sections 3 and 4 the proposed scoring function is employed to address two very different vision problems. In Section 3, the objective is to temporally segment a continuous surveillance video according to the activities taking place in the scene. As ex-



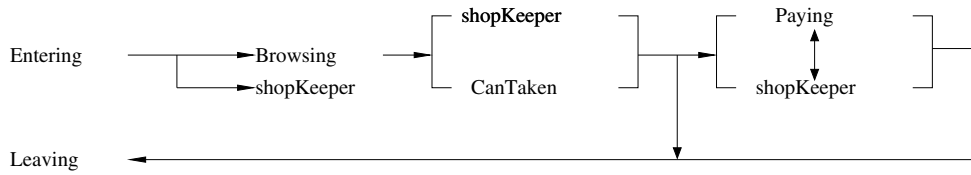
(a) Examples of detected and classified events in a shopping scene



(b) CL-AIC (15 inter-links)

(c) BIC & ICL (13 inter-links)

(d) AIC (19 inter-links)



(e)

Fig. 10. Discovering causal relationships among visual events in a shopping scene. In (a), events belonging to five event classes *canTaken*, *entering/leaving*, *shopkeeper*, *browsing* and *paying* are indicated with bounding boxes in magenta, red, blue, green and cyan, respectively. (b)–(d) Topologies of DML-HMMs learned using different scoring functions. (e) The expected causal and temporal structure of the shopping activity. (For interpretation of color mentioned in this figure the reader is referred to the web version of the article.)

plained earlier, the number of activities that are visually identifiable and can be segmented from a video is determined by the nature of the camera view and the image resolution (i.e., the distance from the camera to the objects of interest). In Section 4, the objective is to discover the causal relationship among the detected events. Obviously, if different number of events were detected (e.g., using different features), the set of causal relationships will be different as well. But, what we are interested in is, for a given set of events, whether the correct causal relationship can be recovered using different graphical model structure scoring functions. So, the focus here is not about what is the optimal number of events one should detect. It is about whether the causal relationship among the actually detected events can be learned correctly.

5. Discussion and conclusion

Our experimental results show that the performance of CL-AIC on learning the structure of a dynamic graphical model with hidden variables is superior to that of existing popular alternatives including BIC, AIC, ICL, and VB. This is especially true when the given data set is noisy and sparse. Similar results were reported in the case of static graphical models in [35]. However, it is interesting to note the difference in the definitions of ‘data sparseness’ in the context of DBNs and in that of static models such as GMMs. The sparseness of a data set is normally measured in comparison to the number of free parameters of a model. The experiments reported in [35] show that a sample size smaller than 5 times of the parameter number should be considered as sparse while our experiments on HMMs and DML-HMMs show that a sample size smaller than 20 times of the number of parameters would qualify for being sparse (see Fig. 2 and 4).

It is worth pointing out that our CL-AIC and the other four scoring functions evaluated in this paper can be classified into two categories based on the underlying principles of derivation. Both AIC and CL-AIC are derived based on information coding theory. Our experimental results suggest that the extra term in the formation of CL-AIC compared to AIC (the second term on the right-hand side of Eq. (4)) is able to rectify the overfitting tendency of AIC given densely distributed data. Different from AIC and CL-AIC, both BIC, ICL, and VB are derived as approximations to Bayesian Model Selection (BMS) [27] which aims to select a model that maximise the marginal likelihood for a data set given the model. Compared to BIC, the approximation made in the derivation of VB is less drastic therefore more accurate resulting in better structure learning. However, it also pays the price for losing generality – VB can only be applied to discrete graphical models whilst the other four can also be implemented for continuous ones such as PCA, ICA, Kalman Filters, and continuous HMMs.

It is also interesting to compare our CL-AIC with a recently proposed structure learning algorithm based on an entropic prior and parameter extinction [13]. The method has been applied to HMM for activity modelling [11]. However, despite being called a structure learning algorithm, it is essentially a parameter learning algorithm aiming to learn the optimal hidden state transition probabilities which are part of the model parameters rather than model structure. With the addition of state trimming, it is possible for the method to learn the optimal number of hidden states for a HMM. Nevertheless, as pointed out in [13], the algorithm tend to keep superfluous gating states in the learned structure. It therefore not suitable for determining the state number for a HMM. Moreover, it cannot be applied for more general dynamical graphical models such as DML-HMMs.

It is critical to distinguish the 'expected model' (the model learned from data) and the true underlying model when studying a real-world model selection problem. In particular, one must note that the expected model and the real (unknown) model underpinning the original data are not necessarily the same. We consider that one of the main causes of this discrepancy is that when solving a real-world data-modelling problem, for computational tractability, one must first perform feature extraction and data pre-processing before a model can be built. However, if those features extracted are not a good representation of the original data, the estimated model, even if it corresponds well to the feature representation, is not a good model for explaining the original data. To make sure that a model does indeed capture accurately the underlying nature of the original data, the choice of feature selection and representation is crucial. We have done exactly that in the two real-data experiments presented in this paper. Specifically, the discrete event based video trajectory representation for the video content analysis has been shown to be more informative and robust to noise compared to conventional representations in our previous work [34]. We have also demonstrated that the event based activity representation is superior to other commonly used representations in revealing the causal/temporal relationships in [19]. We therefore argue with strong justification that for the experiments carried out in this work, the 'expected' model does go some way to explain the true underlying nature of the original problem rather than merely a good explanation of the wanted results.

In conclusion, we have proposed a novel scoring function, CL-AIC for selecting the optimal structure of DBNs. The effectiveness of CL-AIC has been demonstrated on solving the challenging problem of video content analysis. In particular, it is evident from our results that CL-AIC is superior to the commonly used scoring functions including BIC, AIC, ICL, and VB, especially give sparse and noisy visual data. It is worth pointing out that either model parameter estimation or computing the posterior probabilities will become computational intractable for learning the structure of a large DBN. To address this problem, approximation methods such as stochastic simulation [22] or variational approximations [18] need to be adopted. Our ongoing work includes further investigation on how the performance of different scoring functions will be affected by the use of approximation methods.

Appendix. Computing CL-AIC for a DML-HMM

CL-AIC for a DML-HMM is formulated as follows:

$$\text{CL-AIC}(K) = -\log \left(\sum_S p(\mathcal{Y}|S, \hat{\theta}_K) p(S|\hat{\theta}_K) \right) - \log p(S = \hat{\mathcal{Z}}|\mathcal{Y}, \hat{\theta}_K) + C_K. \quad (12)$$

To compute the first and second terms on the right-hand side of the above equation, the ML (maximum likelihood) estimation of the model parameters, $\hat{\theta}_K$ has to be obtained first. This is achieved through an extended dynamic programming algorithm similar to the forward-backward algorithm for a standard HMM.

Let us consider a DML-HMM with C temporal processes and one hidden variable and one observation variable, respectively, for each temporal process at each time instance. We thus have $N_h = N_o = C$ where N_h and N_o are the number of hidden state variables and observation variables at each time instance, respectively. It is assumed that all the hidden state variables are discrete and all the observation variables are continuous whose probability density functions are Gaussian with respect to each state of their parent hidden state variables. The parameter space thus consists of the following components:

1. The initial state distribution $\pi = \{\pi_{i^{(c)}}\}$ where $\pi_{i^{(c)}} = P(S_1^{(c)} = q_{i^{(c)}})$, $1 \leq i^{(c)} \leq N^{(c)}$ and $1 \leq c \leq C$.
2. The state transition probability distribution $\mathbf{A} = \{a_{Pa(j^{(c)})j^{(c)}}\}$, where $a_{Pa(j^{(c)})j^{(c)}} = P(S_{t+1}^{(c)} = q_{j^{(c)}} | Pa(S_{t+1}^{(c)}) = q_{Pa(j^{(c)})})$, $Pa(S_{t+1}^{(c)})$ are the hidden variables at time t on which $S_{t+1}^{(c)}$ is conditionally dependent, $Pa(j^{(c)})$ are subscripts of those discrete values that $Pa(S_{t+1}^{(c)})$ can assume, $1 \leq j^{(c)} \leq N^{(c)}$ and $1 \leq c \leq C$.
3. The observation probability distribution $\mathbf{B} = \{b_{i^{(c)}}(\mathcal{Y}_t^{(c)})\}$, where $b_{i^{(c)}}(\mathcal{Y}_t^{(c)}) = \mathcal{N}(\mathcal{Y}_t^{(c)}; \mu_{i^{(c)}}, \mathbf{U}_{i^{(c)}})$, $\mu_{i^{(c)}}$ and $\mathbf{U}_{i^{(c)}}$ are the mean vector and covariance matrix of the normal (Gaussian) distribution with respect to $S_t^{(c)} = q_{i^{(c)}}$, $1 \leq i^{(c)} \leq N^{(c)}$ and $1 \leq c \leq C$.

Given an observation sequence \mathcal{Y} and a model structure λ , we need to determine the model parameters $\theta(\lambda) = \{\mathbf{A}, \mathbf{B}, \pi\}$ that maximise the probability of the observation sequence given the model structure $P(\mathcal{Y}|\lambda, \theta(\lambda))$. There is no analytical solution to determine the optimal parameters given a finite observation sequence. However, the parameters can be estimated iteratively using dynamic programming. Let us first define the following variables:

- The forward variable $\alpha_t(i^{(1)}, \dots, i^{(C)}) = P(\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_t, S_t^{(1)} = q_{i^{(1)}}, \dots, S_t^{(C)} = q_{i^{(C)}} | \lambda, \theta(\lambda))$, i.e., the probability of the partial observation sequence until time t and states for $S_t^{(1)}, \dots, S_t^{(C)}$, given the model λ and $\theta(\lambda)$:

$$\alpha_t(i^{(1)}, \dots, i^{(C)}) = \begin{cases} \prod_{c=1}^C \pi_{i^{(c)}} b_{i^{(c)}}(\mathcal{Y}_t^{(c)}) & \text{if } t = 1, \\ \prod_{c=1}^C \left(\left(\sum_{j^{(1)}, \dots, j^{(C)}} \alpha_{t-1}(j^{(1)}, \dots, j^{(C)}) a_{Pa(j^{(c)})j^{(c)}} \right) b_{i^{(c)}}(\mathcal{Y}_t^{(c)}) \right) & \text{if } 1 < t \leq T. \end{cases}$$

- The backward variable $\beta_t(i^{(1)}, \dots, i^{(C)}) = P(\mathcal{Y}_t, \dots, \mathcal{Y}_T, S_t^{(1)} = q_{i^{(1)}}, \dots, S_t^{(C)} = q_{i^{(C)}} | \lambda, \theta(\lambda))$, i.e., the probability of the partial observation sequence from $t+1$ to T , given the states for $S_t^{(1)}, \dots, S_t^{(C)}$ and the model λ and $\theta(\lambda)$:

$$\beta_t(i^{(1)}, \dots, i^{(C)}) = \begin{cases} 1 & \text{if } t = T, \\ \sum_{j^{(1)}, \dots, j^{(C)}} \left(\prod_{c=1}^C a_{Pa(j^{(c)})j^{(c)}} b_{j^{(c)}}(\mathcal{Y}_{t+1}^{(c)}) \right) \beta_{t+1}(j^{(1)}, \dots, j^{(C)}) & \text{if } 1 \leq t < T. \end{cases}$$

- $\xi_t(i^{(1)}, \dots, i^{(C)}, j^{(1)}, \dots, j^{(C)}) = P(S_t^{(1)} = q_{i^{(1)}}, \dots, S_t^{(C)} = q_{i^{(C)}} | \lambda, \theta(\lambda), S_{t+1}^{(1)} = q_{j^{(1)}}, \dots, S_{t+1}^{(C)} = q_{j^{(C)}})$, i.e., the probability of being at certain states at time t and $t+1$, given the model and observation sequence:

$$\xi_t(i^{(1)}, \dots, i^{(C)}, j^{(1)}, \dots, j^{(C)}) = \frac{\beta_{t+1}(j^{(1)}, \dots, j^{(C)}) \prod_{c=1}^C \alpha_t(i^{(1)}, \dots, i^{(C)}) a_{Pa(j^{(c)})j^{(c)}} b_{j^{(c)}}(\mathcal{Y}_{t+1}^{(c)})}{P(\mathcal{Y}|\lambda, \theta(\lambda))},$$

where $P(\mathcal{Y}|\lambda, \theta(\lambda))$ can be computed using the forward and backward variables:

$$P(\mathcal{Y}|\lambda, \theta(\lambda)) = \sum_{i^{(1)}, \dots, i^{(C)}} \alpha_t(i^{(1)}, \dots, i^{(C)}) \beta_t(i^{(1)}, \dots, i^{(C)}). \quad (13)$$

- $\gamma_t(i^{(1)}, \dots, i^{(C)}) = P(S_t^{(1)} = q_{i^{(1)}}, \dots, S_t^{(C)} = q_{i^{(C)}} | \mathcal{Y}, \lambda, \theta(\lambda))$, i.e., the probability of being at certain states at time t , given the model and observation sequence:

$$\gamma_t(i^{(1)}, \dots, i^{(C)}) = \frac{\alpha_t(i^{(1)}, \dots, i^{(C)}) \beta_t(i^{(1)}, \dots, i^{(C)})}{\sum_{i^{(1)}, \dots, i^{(C)}} \alpha_t(i^{(1)}, \dots, i^{(C)}) \beta_t(i^{(1)}, \dots, i^{(C)})}. \quad (14)$$

Denote the current parameter estimates as $\theta(\lambda) = \{\mathbf{A}, \mathbf{B}, \pi\}$, the re-estimated parameters $\bar{\theta}(\lambda) = \{\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\pi}\}$ can be computed using the following re-estimation formula: Denote the current parameter estimates as $\theta(\lambda) = \{\mathbf{A}, \mathbf{B}, \pi\}$, the re-estimated parameters $\bar{\theta}(\lambda) = \{\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\pi}\}$ can be computed using the following re-estimation formula:

$$\bar{\pi}_{i^{(c)}} = \sum_{i^{(1)}, \dots, i^{(c-1)}, j^{(c+1)}, \dots, j^{(C)}} \gamma_1(i^{(1)}, \dots, i^{(C)}), \quad (15)$$

$$\bar{a}_{Pa(j^{(c)})|j^{(c)}} = \frac{\sum_{t=1}^{T-1} \sum_{i^{(1)}, \dots, i^{(c-1)}, j^{(c+1)}, \dots, j^{(C)}, i^{(c)} \neq Pa(j^{(c)})} \zeta_t(i^{(1)}, \dots, i^{(C)}, j^{(1)}, \dots, j^{(C)})}{\sum_{t=1}^T \sum_{i^{(c')} \neq Pa(j^{(c)})} \gamma_t(i^{(1)}, \dots, i^{(C)})}, \quad (16)$$

$$\bar{\mu}_{i^{(c)}} = \frac{\sum_{t=1}^T \left(\sum_{i^{(1)}, \dots, i^{(c-1)}, j^{(c+1)}, \dots, j^{(C)}} \gamma_t(i^{(1)}, \dots, i^{(C)}) \right) \mathcal{Y}_t^{(c)}}{\sum_{t=1}^T \sum_{i^{(1)}, \dots, i^{(c-1)}, j^{(c+1)}, \dots, j^{(C)}} \gamma_t(i^{(1)}, \dots, i^{(C)})}, \quad (17)$$

$$\bar{\mathbf{U}}_{i^{(c)}} = \frac{\sum_{t=1}^T \left(\sum_{i^{(1)}, \dots, i^{(c-1)}, j^{(c+1)}, \dots, j^{(C)}} \gamma_t(i^{(1)}, \dots, i^{(C)}) \right) (\mathcal{Y}_t^{(c)} - \bar{\mu}_{i^{(c)}})(\mathcal{Y}_t^{(c)} - \bar{\mu}_{i^{(c)}})^T}{\sum_{t=1}^T \sum_{i^{(1)}, \dots, i^{(c-1)}, j^{(c+1)}, \dots, j^{(C)}} \gamma_t(i^{(1)}, \dots, i^{(C)})}. \quad (18)$$

If we iteratively use $\bar{\theta}(\lambda)$ to replace $\theta(\lambda)$ and repeat the re-estimation calculation until some limiting point is reached, the final result is a maximum likelihood estimate of parameters $\hat{\theta}(\lambda)$.

Now given the learned parameters $\hat{\theta}(\lambda)$, the first term on the right-hand side of Eq. (12) can be computed using Eq. (13). The remaining problem is to compute the second term which involves the estimation of the most probable hidden state sequence. We formulate an extended Viterbi algorithm to infer the sequence given the observation and learned parameters. Let us first define the following variables:

- $\delta_t(i^{(1)}, \dots, i^{(C)}) = \max_{\{S_t^{(c)}\}, \dots, \{S_{t-1}^{(c)}\}, S_t^{(1)} = q_{i^{(1)}}, \dots, S_t^{(C)} = q_{i^{(C)}}, \mathcal{Y}_1, \dots, \mathcal{Y}_t | \lambda, \theta(\lambda)}$, where $\{S_t^{(c)}\} = \{S_t^{(1)}, \dots, S_t^{(C)}\}$. $\delta_t(i^{(1)}, \dots, i^{(C)})$ is the highest probability of the partial observation sequence until time t and a sequences of hidden states from time 1 to time t , given the model λ and $\theta(\lambda)$.
- $\varphi_t(i^{(1)}, \dots, i^{(C)})$, which is an array used to store the best state sequence. The best hidden states at time t for the C hidden state variables are denoted as $\{S_t^{(c)}\}$.

The extended Viterbi algorithm has the following steps:

1. Initialisation:

$$\delta_1(i^{(1)}, \dots, i^{(C)}) = \prod_{c=1}^C \pi_{i^{(c)}} b_{i^{(c)}}(\mathcal{Y}_1^{(c)}),$$

$$\varphi_1(i^{(1)}, \dots, i^{(C)}) = \mathbf{0}.$$

2. Recursion:

$$\delta_t(i^{(1)}, \dots, i^{(C)}) = \max_{j^{(1)}, \dots, j^{(C)}} \left\{ \prod_{c=1}^C \delta_{t-1}(j^{(1)}, \dots, j^{(C)}) a_{Pa(i^{(c)})|j^{(c)}} b_{i^{(c)}}(\mathcal{Y}_t^{(c)}) \right\},$$

$$\varphi_t(i^{(1)}, \dots, i^{(C)}) = \arg \max_{j^{(1)}, \dots, j^{(C)}} \left\{ \prod_{c=1}^C \delta_{t-1}(j^{(1)}, \dots, j^{(C)}) a_{Pa(i^{(c)})|j^{(c)}} b_{i^{(c)}}(\mathcal{Y}_t^{(c)}) \right\},$$

where $1 < t \leq T$.

3. Termination:

$$\{S_T^{(c)}\} = \arg \max_{i^{(1)}, \dots, i^{(C)}} \delta_T(i^{(1)}, \dots, i^{(C)}).$$

4. Best state sequence backtracking:

$$\{S_t^{(c)}\} = \varphi_{t+1}(\{S_{t+1}^{(c)}\}),$$

where $t = T - 1, T - 2, \dots, 1$.

Given the most probable hidden states sequence $\hat{\mathcal{Z}} = \{S_t^{(c)}\}$, the second term of Eq. (12) is computed as

$$\log p(S = \hat{\mathcal{Z}} | \mathcal{Y}, \hat{\theta}_K) = \sum_{t=1}^T \log \gamma_t, \quad (19)$$

where γ_t , computed using Eq. (14), is the probability of being at the most probable hidden states sequence at time t given the model and observation sequence.

References

- [1] H. Akaike, Information theory and an extension of the maximum likelihood principle, in: Proceedings of the 2nd International Symposium on Information Theory, 1973, pp. 267–281.
- [2] N. Babaguchi, Y. Kawai, T. Kitahashi, Event based indexing of broadcasting sports video by intermodal collaboration, IEEE Transactions on Multimedia 4 (1) (2002) 68–75.
- [3] L.E. Baum, T. Petrie, Statistical inference for probabilistic functions of finite state Markov chains, Annals of Mathematical Statistics 37 (1966) 1554–1563.
- [4] M.J. Beal, Z. Ghahramani, Variational bayesian learning of directed graphical models with hidden variables, Bayesian Analysis 1 (2006) 1–44.
- [5] M. Beal, Z. Ghahramani, The variational Bayesian EM algorithm for incomplete data: with application to scoring graphical model structures, Bayesian Statistics 7 (2003).
- [6] C. Biernacki, G. Celeux, G. Govaert, Assessing a mixture model for clustering with the integrated completed likelihood, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (7) (2000) 719–725.
- [7] C. Biernacki, G. Govaert, Using the classification likelihood to choose the number of clusters, Computing Science and Statistics 29 (1997) 451–457.
- [8] J. Boreczky, L. Wilcox, A hidden markov model framework for video segmentation using audio and image features, in: International Conference on Acoustics, Speech, and Signal Processing, 1998, pp. 3741–3744.
- [9] C. Boutilier, N. Friedman, M. Goldszmidt, and D. Koller, Context-specific independence in bayesian networks, in: Uncertainty in AI, 1996, pp. 115–123.
- [10] R. Boys, D. Henderson, A bayesian approach to dna sequence segmentation, Biometrics 60 (2004) 573–588.
- [11] M. Brand, V. Kettmaker, Discovery and segmentation of activities in video, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (8) (2000) 844–851.
- [12] M. Brand, N. Oliver, A. Pentland, Coupled hidden markov models for complex action recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, Puerto Rico, 1996, pp. 994–999.
- [13] M. Brand, Structure discovery in conditional probability models via an entropic prior and parameter extinction, Neural Computation 11 (5) (1999) 1155–1182.
- [14] F. Brugnara, D. Falavigna, M. Omologo, Automatic segmentation and labeling of speech based on hidden markov models, Speech Communication 12 (1993) 357–370.
- [15] T. Duong, H. Bui, D. Phung, and S. Venkatesh, Activity recognition and abnormality detection with the switching hidden semi-markov model, in: IEEE Conference on Computer Vision and Pattern Recognition, 2005, pp. 838–845.
- [16] G.D. Forney, The Viterbi algorithm, Proceedings of the IEEE 61 (1973) 268–278.
- [17] N. Friedman, K. Murphy, S. Russell, Learning the structure of dynamic probabilistic networks, in: Uncertainty in AI, 1998, pp. 139–147.
- [18] Z. Ghahramani, Learning dynamic bayesian networks, in: Adaptive Processing of Sequences and Data Structures, Lecture Notes in AI, 1998, pp. 168–197.
- [19] S. Gong and T. Xiang, Recognition of group activities using dynamic probabilistic networks, in: IEEE International Conference on Computer Vision, 2003, pp. 742–749.
- [20] D. Heckerman, D. Geiger, D. Chickering, Learning bayesian networks: the combination of knowledge and statistical data, Machine Learning 20 (1995) 197–243.
- [21] N. Johnson, A. Galata, D. Hogg, The acquisition and use of interaction behaviour models, in: IEEE Conference on Computer Vision and Pattern Recognition, 1998, pp. 866–871.
- [22] K. Kanazawa, D. Koller, S. Russell, Stochastic simulation algorithms for dynamic probabilistic networks, in: Uncertainty in AI, 1995, pp. 346–351.
- [23] R. Kass, A. Raftery, Bayes factors, Journal of the American Statistical Association 90 (1995) 377–395.
- [24] S. Kullback, Information Theory and Statistics, Dover, New York, 1968.
- [25] L.R. Rabiner, A tutorial on hidden Markov models and selected applications in speech recognition, Proceedings of the IEEE 77 (2) (1989) 257–286.
- [26] N. Oliver, B. Rosario, A. Pentland, A Bayesian computer vision system for modelling human interactions, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (8) (2000) 831–843.
- [27] A. Raftery, Bayes model selection in social research, Sociological Methodology 90 (1995) 181–196.
- [28] J. Rissanen, Stochastic Complexity in Statistical Inquiry, World Scientific, Singapore, 1989.
- [29] G. Schwarz, Estimating the dimension of a model, Annals of Statistics 6 (1978) 461–464.
- [30] R. Shibata, Selection of the order of an autoregressive model by Akaike's Information Criterion, Biometrika 63 (1976) 117–126.
- [31] C. Stauffer, W. Grimson, Learning patterns of activity using real-time tracking, IEEE Transactions on Pattern Analysis and Machine Intelligence 22 (8) (2000) 747–758.

- [32] C. Vogler, D. Metaxas, A framework for recognizing the simultaneous aspects of American sign language, *Computer Vision and Image Understanding* 81 (2001) 358–384.
- [33] T. Xiang, S. Gong, D. Parkinson, Autonomous visual events detection and classification without explicit object-centred segmentation and tracking, in: *British Machine Vision Conference*, 2002, pp. 233–242.
- [34] T. Xiang, S. Gong, Activity based video content trajectory representation and segmentation, in: *British Machine Vision Conference*, 2004, pp. 177–186.
- [35] T. Xiang, S. Gong, Visual learning given sparse data of unknown complexity, in: *IEEE International Conference on Computer Vision*, 2005, pp. 701–708.
- [36] T. Xiang, S. Gong, Optimal dynamic graphs for video content analysis, in: *British Machine Vision Conference*, 2006.