

AUTOMATICALLY LEARNING STRUCTURAL UNITS IN EDUCATIONAL VIDEOS WITH THE HIERARCHICAL HIDDEN MARKOV MODELS

Dinh Q. Phung, Svetha Venkatesh

School of Computing
Curtin University of Technology
GPO Box U1987, Perth, 6845, Western Australia
{phungquo,svetha}@computing.edu.au

Hung H. Bui

Artificial Intelligence Center
SRI International, 333 Ravenswood Ave
Menlo Park, CA 94025, USA
bui@ai.sri.com

ABSTRACT

In this paper we present a coherent approach using the hierarchical HMM with shared structures to extract the structural units that form the building blocks of an education/training video. Rather than using hand-crafted approaches to define the structural units, we use the data from nine training videos to learn the parameters of the HHMM, and thus naturally extract the hierarchy. We then study this hierarchy and examine the nature of the structure at different levels of abstraction. Since the observable is continuous, we also show how to extend the parameter learning in the HHMM to deal with continuous observations.

1. INTRODUCTION

The indexing of multimedia content is an important issue in the management of vast amount of data. This involves the ability to automatically categorise the data at different levels of semantic abstraction by combining the cues from all modalities. Whilst this categorisation can be achieved to some extent by a bottom-up approach that uses the visual and aural media explicitly, it is our contention that higher levels of abstraction can be extracted if a data driven bottom up approach is combined with top-down contextual domain knowledge.

In this paper we examine a category of video, the education and training videos, whose objective is to educate and train. This particular video genre is important for e-learning and is often a simpler category of video as the aim of the video is unambiguous, and the techniques for educating and training fall in a somewhat well contained set, as opposed to a much larger set of techniques for handling aesthetics in motion pictures. The narrative in such video is also simpler, often linear, and mostly developed in a hierarchic way to elucidate concepts in greater levels of detail. Deductive approaches will move from the general to the particular, whilst inductive approaches will move from the detail to the general. The set of conventions and rules that shape these videos is termed production grammar, and our earlier work [1] has explored the structural elements. In this work we present a coherent approach using the hierarchical HMM to extract the structural units that form the building blocks of an education/training video. Rather than using the hand-crafted approach of our previous work [1] to define the structural units, we use the data to learn the parameters of a HHMM, and thus naturally extract the hierarchy. We then study this hierarchy and examine the nature of the structural units at different levels of abstraction.

Discovering structure in videos is a current active research area. The central task is to efficiently index data into semantic units, possibly at different level of abstractions, to simplify the process of navigation. There have been several systems proposed for specific video genres. Partition and classification of broadcast videos into meaningful sections has attracted significant attention. For instance, in [2], low-level features are combined with the concept of shot syntax to identify and label different narrative structures such as anchor shots, voice-over segments and interview sections found in news programs. Research into the domain of lecture videos be found in [3]. In their work, visual events are detected from the visual stream and then incorporated with audio information in a probabilistic framework to detect topic transitions. The domain of entertainment film has also been studied lately. Wang *et al.* [4], for example, attempt to detect *scenes* in film using the similarity in visual information and further improve the results with guidance from cinematic grammar.

In all videos there is often a common and natural hierarchy in the content. For example, surveillance videos map to different regions; news reports are organised into layers of details; a film has episodes, story units, scenes then shots, and a training video has several sub-topics which in turn form a main topic. Being able to model the hierarchic nature of a video offers ways to index the content in a meaningful hierarchy of semantic units. The hierarchical HMM has emerged naturally as a candidate to model this problem. The HHMM is a powerful stochastic model, first introduced in [5], in which the HHMM is viewed as a form of probabilistic context free grammar (PCFG), and the inference algorithm and parameter learning procedures are constructed based on the inside-outside algorithm. In [6], the HHMM is converted to a DBN, and a general DBN inference algorithm is applied to achieve complexity linear in time T , but exponential in the depth D of the HHMM. The same approach is applied in [7], ie: the HHMM is 'flattened' into regular HMM with a very large state space for inference purpose. Their work [7] aims to detect structures of soccer videos in an unsupervised manner. The model selection is first carried out using the MCMC to determine the structure parameters for the model, followed by a feature selection procedure. Finally, the HHMM is used to detect two semantic concepts, namely *play* and *break* in soccer videos. As there is little hierarchy at this level, the power of the hierarchic probabilistic model is not used. The HHMM is also applied in other domains other than multimedia such as in hand-written recognition [5], robot navigation, behaviour recognition and information retrieval.

The remainder of the paper is organised as follows. We briefly

discuss the HHMM with shared structures in Sec. 2. Next, Sec. 3 presents a brief analysis of the narrative structure in educational videos. This is followed by the experimental results in Sec. 4. Finally, the conclusion is provided in Sec. 5

2. THE HIERARCHICAL HMM WITH SHARED STRUCTURES

The discrete HHMM was originally proposed in [5] and its extension to accommodate shared structured has been addressed in our previous work [8]. We refer readers to [8] for a complete treatment. Here, we refocus our attention to elucidate the idea of the shared structures and briefly discuss the DBN representation and the EM procedure for parameter estimation. We then present the case when the emission probability is modeled as a mixture of Gaussians.

A hierarchical HMM extends the traditional HMM to allow each state itself to be a HMM. Formally, a HHMM is defined by a topological structure Γ and a set of parameter θ attached to the topology. The depth D —defines the number of layers in the hierarchy—and the *number of states* available at each level Q^d , for $d = 1, \dots, D$, are specified by Γ . Fig. 2 shows an example of a topology of depth 3. Level 1 is the root level and is always fixed to have only a single state, and only the lowest level D , termed as *production level*, emits observation. Furthermore, the topological structure reveals the ‘parent-children’ relationship of states between two consecutive levels¹. A state p at level d is assigned to a *set of children*, $\text{ch}(p)$, at level $d + 1$. A state i at level $d + 1$ therefore might multiply inherit, or be *shared*, from more than one parent at level d . The general form of DBN representation² is shown in Fig. 1(a).

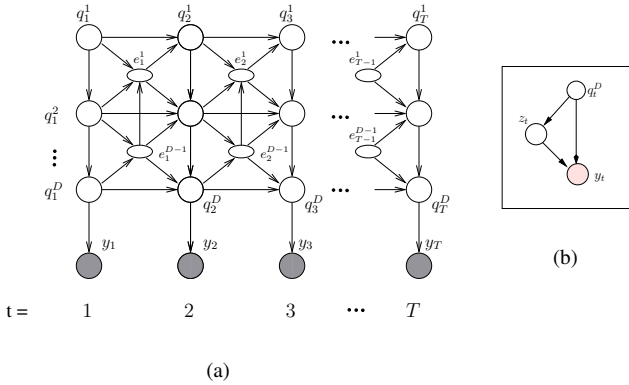


Fig. 1. (a) DBN representation for the discrete HHMM; (b) Added mixture component z_t at level D .

Given a topology Γ , the parameter $\theta \triangleq \{\pi, A, \mu, \Sigma\}$ of the HHMM is specified as follows.. For each level d (except D), $p \in Q^d$, and $i, j \in \text{ch}(p)$ we define: $\pi_{i,j}^{d,p}$ is the initial probability of the child i given the parent is p at level d ; $A_{i,j}^{d,p}$ is the transition probability from child i ; and $A_{i,\text{end}}^{d,p}$ is the probability that state p terminates at level d given its current child is i . We require that $\sum_i \pi_{i,j}^{d,p} = 1$, $\sum_j A_{i,j}^{d,p} = 1$, and $A_{i,\text{end}}^{d,p} \leq 1$, which later become Lagrangian constraints used in the maximisation step in the EM procedure. Finally, at level D an observation probability matrix

¹Note that the original HHMM in [5] assumes that a state has a only a single parent and therefore the topology reduces strictly to a tree.

²The DBN structure for HHMM was originally presented in [6]

B is specified in the case of discrete HHMM, or a set of $\{\mu, \Sigma\}$ in case the observation values are continuous and modeled as a Gaussian mixture model (GMM).

Given an observed data set \mathcal{D} and some initial parameters, the EM algorithm iteratively re-estimates a new parameter $\hat{\theta}$, hill climbing in the parameter space which is guaranteed to converge to a local maxima. As shown in [8], doing EM parameter estimation reduces to first calculating the expected sufficient statistics (ESS) $\bar{\tau} = E_{\nu \setminus \mathcal{D}} \tau$, and then setting the re-estimated parameter $\hat{\theta}$ to the normalized value of $\bar{\tau}$. The ESS for parameter $\{A_{i,j}^{d,p}\}$, for example, is calculated as:

$$\bar{\tau}(A)_{i,j}^{d,p} = E_{\nu \setminus \mathcal{D}} \tau(A)_{i,j}^{d,p} = \sum_{t=1}^{T-1} \xi_t^{d,p}(i, j) / \Pr(\mathcal{D}) \quad (1)$$

where the auxiliary variable $\xi_t^{d,p}(i, j)$ is defined as the transition probability $\Pr(q_{t+1}^{d+1} = j, q_t^{d+1} = i, q_{t+1}^d = p, e_t^{d:d+1} = 01, \mathcal{D})$, which is diagrammatically visualised as $\left[\begin{smallmatrix} i \\ \text{---} \\ \text{---} \\ \text{---} \\ j \end{smallmatrix} \right]$. We refer readers to [8] for complete details on the computation of auxiliary variables and other expected sufficient statistics. The rest of this section will discuss the case when the emission probability is modeled as a GMM.

The general method when modeling GMM as the observation probability in the hierarchical HMM is similar to that of the regular HMM. The DBN structure at level D is modified as in Fig 1(b), where a mixture variable z_t is added. For simplicity, thereafter in this section we will drop the index D . Let M be the number of mixtures and N be the number of states at level D . The observation matrix B in the discrete case is replaced by the mixing weight matrix $\{\lambda_{mi}\}$ and a set of means and covariance matrices $\{\mu_{mi}, \Sigma_{mi}\}$ for $i = 1, \dots, N$ and $m = 1, \dots, M$. Given observed data \mathcal{D} , following a standard EM procedure, the expected complete log-likelihood $\langle \ell(\theta; \mathcal{D}) \rangle$ in the E-step is given as (discarding terms irrelevant to z_t and y_t):

$$\sum_{\substack{1 \leq i \leq N \\ 1 \leq m \leq M}} \left[\sum_{t=1}^T \langle I_{m,i}^{z_t, q_t} \rangle \log \mathcal{N}(y_t | \mu_{mi}, \Sigma_{mi}) + \sum_{t=1}^T \langle I_{m,i}^{z_t, q_t} \rangle \log \lambda_{mi} \right]$$

where $\mathcal{N}(\cdot)$ is the multivariate Gaussian density function, $\langle \cdot \rangle$ is the expectation operator, and $I_{m,i}^{x,y}$ is the identity function which takes 1 if $\{x = m \cup y = i\}$, and 0 otherwise. It is straightforward to compute the expectation of the identity function:

$$\begin{aligned} \langle I_{m,i}^{z_t, q_t} \rangle &= \Pr(z_t = m, q_t = i | \mathcal{D}) = \Pr(z_t = m | q_t = i, y_t) \Pr(q_t = i | \mathcal{D}) \\ &= \frac{\lambda_{mi} \mathcal{N}(y_t | \mu_{mi}, \Sigma_{mi})}{\sum_{m=1}^M \lambda_{mi} \mathcal{N}(y_t | \mu_{mi}, \Sigma_{mi})} \times \frac{\gamma_t^D(i)}{\Pr(\mathcal{D})} \end{aligned}$$

where³ the auxiliary variable $\gamma_t^D(i)$ defined as the probability $\Pr(q_t^D = i, \mathcal{D})$ and can be computed directly from horizontal transition probability $\xi_t^{d,p}(i, j)$ and vertical transition probability $\chi_t^{d,p}(i)$ (defined [8]). Finally, in M-step we maximise the expected complete log-likelihood $\langle \ell(\theta; \mathcal{D}) \rangle$ with respect to λ_{im} and μ_{mi}, Σ_{mi} . Introducing Lagrange multipliers for λ_{mi} and setting derivatives to zero for μ_{mi}, Σ_{mi} , the set of re-estimated parameters is given

³we put back hierarchic index D for clarity here

as:

$$\hat{\epsilon}_{mi} = \frac{\sum_{t=1}^T \langle I_{m,i}^{z_t,qt} \rangle}{\sum_{m=1}^M \sum_{t=1}^T \langle I_{m,i}^{z_t,qt} \rangle}, \quad \hat{\mu}_{mi} = \frac{\sum_{t=1}^T \langle I_{m,i}^{z_t,qt} \rangle y_t}{\sum_{t=1}^T \langle I_{m,i}^{z_t,qt} \rangle}$$

$$\hat{\Sigma}_{mi} = \frac{\sum_{t=1}^T \langle I_{m,i}^{z_t,qt} \rangle (y_t - \mu_{mi})(y_t - \mu_{mi})^T}{\sum_{t=1}^T \langle I_{m,i}^{z_t,qt} \rangle}$$

When multiple (iid) observation sequences are given, the set of above equations can be adjusted by simply adding a summation over the number of sequences. This corresponds to ‘counting’ over all sequences.

3. NARRATIVE STRUCTURAL UNITS

The structure of an education/training video is fashioned by the way in which training material is presented. In almost all cases, there is a narrator that takes us on the journey of learning through the subject. Thus, the narrator will reappear through the video, guiding us along the way. To make the points, the narrator uses a bag of tricks: text, footage with voice-over, or discussion with people or interviews. Generally a shot in this video genre can be classified into one of following three categories. *On-screen* narration refers to sections of narration, both direct and assisted, in which the narrator directly speaks to the audience. *Voice-over* refers to sections where the narrator will maintain the narrative by using voice-over. *Linkage* sections refers to sections of informational footage without the narrator and we divide this into two subgroups: *expressive linkage* and *functional linkage*, where the former is used to dramatise a situation. Readers are referred to [1] for precise definitions and examples of these structural units. For now, it is sufficient for us to augment these four main categories as a set of high level semantics. In a generative process, one can imagine that they correspond to four parent states that consists of a number children states producing what we see in the video such as the red colour, the caption texts and so on. We use this knowledge to construct a topology for the experiment presented next.

4. EXPERIMENTAL RESULTS

In this experiment, we present initial results on applying the HHMM to automatically map semantic concepts to the model states at different levels of abstraction. Taking the advantages of the expressiveness of shared structures, we construct a 3-level HHMM with a ‘true’ model in mind, ie: the topology of the HHMM is constructed in such a way that it maps into the hypothesised hierarchy as in Fig 2. That is, for example, at the production level, we are hypothesise that the model can ‘perfectly’ learn to map each state to an elementary feature such as face, text or speech, from which higher levels of descriptions such as *on-screen* narration section can be built.

Given this topology, a 3-level HHMM is subsequently formed. The production level includes 7 states directly attached to the observed features (recall that these are shot-based features). The next level have 4 states, which in our opinion, adequately reflects the number of higher level units in this experiment. Finally, the root covers the entire video. We first randomly initialise the parameter of the model and then use EM to learn a new set of model parameters.

The feature vector used in this experiment includes seven features computed at the shot level. Three of them are from the visual

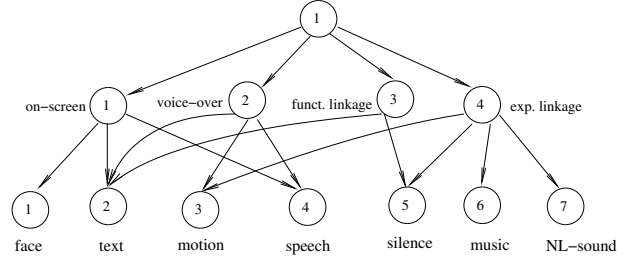


Fig. 2. Hierarchy of connected concepts which is used as topological specification for the corresponding HHMM (links from state 1, 2 (level 2) to state 7 (level 3) are not shown for readability).

stream, namely: face-content-ratio, text-content-ratio, and average motion based on camera pan and tilt estimation; and the other four features are extracted from the sound track including: music-ratio, speech-ratio, silence-ratio, and non-literal sound (NL-) ratio. Computation of these features are detailed in [9]. Note that all of these features are in the range of [0, 1].

Let V be any arbitrary video, and T be the number of shots in V . Feature extraction from V will result in a sequence of observations of length T where each observation is a column vector of seven features. We collect nine videos in total to use for the training purpose with T ranging from 124 to 245. We analyse the learned model and present the results at two levels.

Semantic mapping at the production level

We use the mean matrix at the production level to understand the mapping between the production level states and the feature vector components. Figure 3 shows the visualisation of the means. Thus, for example, state 1 is strongly linked to face and speech, and a little bit of motion. This state corresponds to direct and assistive narration sections. State 2 corresponds to speech and motion, and thus to voice-over sections. State 3 corresponds to NL-sound, and thus to expressive linkage sections and so on. In Table 1, we summarise this result of mapping the production level states to the structural units manually crafted in our earlier work. The mean values learned at the production level is visualised in Fig 3.

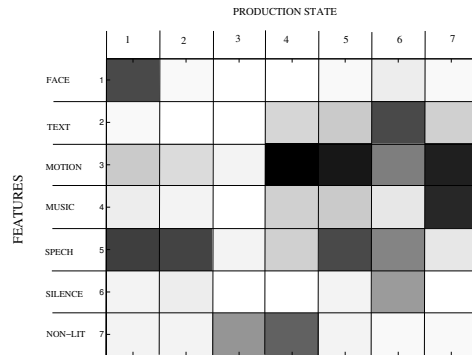


Fig. 3. Gray-scaled visualisation of the mean values of seven features learned at the production level.

Semantic mappings at the upper level

The middle level of topology (Fig. 2) contains 4 states. Combining the information of the ‘parent-child’ relationship specified in topology and results from Table 1, we can interpret the meaning of

