

Segmenting Time Series with a Hybrid Neural Networks - Hidden Markov Model

Abstract

This paper describes work on a hybrid HMM/ANN system for finding patterns in a time series, where a pattern is a function realizable by a recurrent neural network embedded in the state of a hidden Markov model. The most likely (Viterbi) path of the hidden Markov model is used both for re-training the HMM/ANN model and for segmenting the time series into pattern occurrences. The number of patterns is determined from the data by increasing the number of networks as long as the likelihood of the segmentation increases. Preliminary results show that potentially useful patterns can be induced this way.

1 Motivation

An agent embodied in a robot receives information about its interactions with the environment as time series of sensor vectors. As a first step in learning to represent the world, we want the agent to learn to identify segments of the time series that correspond to distinct parts of interactions, and create low-level concepts associated with these segments. For example, figure 1 shows that while approaching an object, the evolution of the $vis-A-x$ sensor, the object's x -coordinate in the robot's visual field, is well predicted by a function $vis-A-x(t+1) = f(trans-vel(t), vis-A-x(t))$, where $trans-vel$ is the robot's translational velocity. When the robot passes the object, f 's prediction is no longer accurate. We would like the agent to induce concepts like approaching or passing an object and recognize them in its future experiences. The example shows that a concept can be represented by a continuous function that predicts future sensor values. We will call such a concept a function pattern. The goal of this work is to induce this kind of patterns. Given a time series of sensor data, the learning task is to find an appropriate set of function patterns and to identify the times at which each pattern occurs, thus producing a segmentation of the time series.

Because they can approximate arbitrarily well continuous functions, we choose artificial neural networks (ANN) to represent the function patterns. There are certain regularities in the sequence of patterns – for example, during an

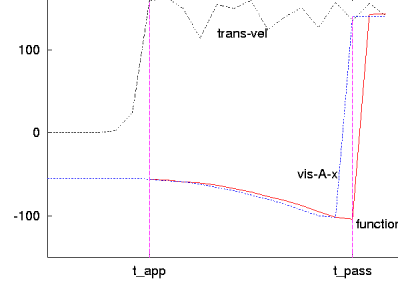


Figure 1. Sensor prediction during a “pass-right” experience; the function in this case is

$$vis-A-x(t+1) = vis-A-x(t) \left(1 - \frac{1}{1-c*trans-vel(t)} \right).$$

object passing experience, a “no object” pattern usually follows an “approach” pattern, and we chose Markov chains to model these regularities. The resulting model is a hidden Markov model (HMM) whose states are neural networks: at each moment the data generating process is assumed to be in one HMM state, and the observed value is assumed to be generated by the neural network embedded in that state. The goal of the induction algorithm is to find the model that produces the maximum likelihood segmentation.

2 The learning framework and algorithms

While motivated by the goal of having a robot agent learn representations, the learning paradigm described here applies to time series generated by other processes as well, so it will be presented in general terms. Given a time series of observed vectors $x(1) \dots x(t) \dots x(T)$ with $x(t) \in R^n$, we assume that it was generated by a set of K processes – or function patterns, each process k being described by a continuous function $f_k = R^n \rightarrow R^d$:

$$\begin{aligned} y(t) &= f_k(x(t), x(t-1), \dots) \\ x_{i_j}(t+1) &= y_j(t) + e_k(t) \quad 1 \leq j \leq d \end{aligned}$$

with $e_k(t)$ a random normal variable representing noise. Variables x_{i_j} , whose probability densities are controlled by

functions f_k , are said to be the output, or predicted, variables. The remaining variables are considered input variables provided by the environment; their probability distribution will be ignored here. The partitioning into input and output variables is considered given. A data point at time t is the pair $\langle x(t), y(t) \rangle$, with $y \in \mathbb{R}^d$ the vector of output variables. The pair $\langle x(t), y(t) \rangle$ will be denoted by $o(t) \in \mathbb{R}^{n+d}$, and the resulting time series by the sequence $O = o(1) \dots o(t) \dots o(T)$. The likelihood of observing $o(t)$ within O , given function f_k , is the value of the multivariate normal density $N(e_k(t), 0, \Sigma^k)$:

$$\rho_k(o(t)) = \frac{1}{(2\pi|\Sigma_k|)^{d/2}} e^{-\frac{e_k(t)^t \Sigma_k^{-1} e_k(t)}{2}} \quad (1)$$

$$e_k(t) = y(t) - f_k(x(t), x(t-1), \dots) \quad (2)$$

where $e_k(t)^t$ is the transpose of vector $e_k(t)$. The parameter Σ^k is a diagonal covariance matrix (i.e. the noise variables are assumed uncorrelated) associated with process k . The process of switching from one function pattern to another is assumed to be described by a stationary first-order Markov chain. This is a strong assumption that we intend to relax in future work, but for now it provides a computationally tractable model. Each function pattern k is a state of the Markov process and has a set $a_k = \{a_{k,l}, 1 \leq l \leq K\}$ of state transition probabilities. An additional parameter $a_{0,k}$ is the probability that k is the initial state. The resulting structure $\lambda = \{s_k = \langle f_k, \Sigma_k, a_{0,k}, a_k \rangle, 1 \leq k \leq K\}$ is a hidden Markov model ([3]) with K states $\{s_k\}$. Given a model λ , a state path $S = s(1) \dots s(T)$ specifying the state of the process at each time, and an observed time series O , the likelihood of O being generated by λ along path S is:

$$L_\lambda(O, S | \lambda) = \prod_{t=1}^T a_{s(t-1), s(t)} * \rho_{s(t)}(o(t))$$

For time series O and model λ , the best segmentation of O is considered to be the one given by the most likely path:

$$V = \arg \max_S L_\lambda(O, S | \lambda)$$

This path is called the Viterbi path because it is computed by the Viterbi algorithm [3].

We assume that the functions f_k can be computed by recurrent neural networks. Because neural networks can approximate arbitrarily well any continuous function, and even some discontinuous functions ([1]), this is a weak (non-restrictive) assumption.

The learning task can now be described as finding a model $\lambda = \{\langle net_k, \Sigma_k, a_k \rangle, 1 \leq k \leq K\}$ that maximizes the likelihood $L_\lambda(O, V_\lambda)$ for a given time series O ; net_k is the neural network of state k . The subscript in V_λ indicates that the segmentation depends on the model λ .

2.1 RECURRENT NEURAL NETWORK ARCHITECTURE AND TRAINING

Expression 2 shows that in the general case the output value of network net_k (which computes function f_k) depends on the entire past sequence of observed variables. In order to implement this variable length memory, recurrent neural networks with the architecture depicted in figure 2 were chosen. The input and output units have identity

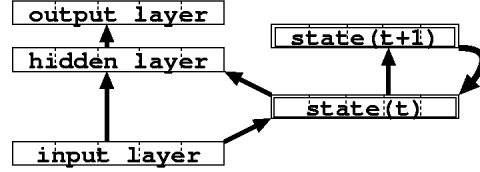


Figure 2. Recurrent network architecture: $s(t+1) = func(s(t), x(t))$ is the network's internal state, encoded in the recurrent units.

transfer functions, the hidden and state units have cosine activation functions. The cosine units can readily construct an approximating function – since they form a function's Fourier representation (see Barron [1]), and unlike the more popular sigmoidal units, they do not saturate during gradient descent. State units can also have identity transfer functions. The network architecture (size, topology, transfer functions) is predetermined, not adapted during training.

Assuming an initial segmentation of the time series, each network is trained to estimate the y values within its assigned segments. The networks are trained with gradient descent on the error surface. The error function is the normalized mean square error, $E = \frac{1}{2l} \sum_{t=1}^l \sum_{j=1}^d \frac{(y_j(t) - z(t))^2}{\sigma_j}$,

where $z(t) = f(x(t), t)$ is the network output, l is the number of points allocated to the network, and the first sum is taken over these points. This particular error function is the one that must be minimized when maximizing the likelihood $L_\lambda(O, V)$.

2.2 HMM induction

Given a segmentation V , and K functions $\{f^k\}$ implemented by neural networks, the model parameters $\{\Sigma_k, a_k\}$ can be estimated by maximizing the likelihood $L_\lambda(O, V)$. The maximum can be found by setting the partial derivatives of the likelihood to 0. The covariance matrix Σ_k is assumed diagonal, so we need to estimate only the variances σ_j^k of the y_j variables in state k . The solutions are the average network errors:

$$\sigma_j^k = \frac{1}{size(T_k)} \sum_{t \in T_k} (y(t) - f_k(x(t), t))^2 \quad (3)$$

where \mathcal{L}_k is the set of (indices of) points allocated to network k in segmentation V . The solutions for the probability transitions are:

$$a_{k,l} = \frac{\#s_k \rightarrow s_l}{\#s_k^-} \quad (4)$$

where $\#s_k \rightarrow s_l$ is the number of transitions from state s_k to state s_l in path V , and $\#s_k^-$ is the number of occurrences of s_k in V (not counting V 's last element).

2.3 Putting it all together: training both the networks and the HMM

The goal is to find the model that maximizes the likelihood of the observed time series along its Viterbi path. Algorithm 1 searches heuristically for a local maximum by adding new network states as long as the likelihood keeps increasing.

Algorithm 1 Estimation of λ for maximizing $L_\lambda(O, V_\lambda)$:

- initialization: start with an arbitrary λ with m states
 - repeat:
 - $\lambda \leftarrow$ the best model with m states (algorithm 2)
 - calculate V_λ
 - poorly predicted points along $V_\lambda \Rightarrow$ create p new networks and randomly assign these points to them (change segmentation), $m \leftarrow m + p$
- until the likelihood stops increasing

A point is considered poorly predicted if its error is greater than a dynamically computed threshold.

Algorithm 2, which finds a (local) best model with a fixed number of states, is related to the expectation maximization (EM) algorithm. It differs in that instead of trying to maximize the model's expected likelihood, it tries to maximize the model's maximum likelihood – the likelihood along the model's Viterbi path.

Algorithm 2 Best model with fixed number of states:

- start with λ and V_λ
 - repeat:
 1. estimate a new model λ^* from V_λ : train the networks, calculate the variances and the transition probabilities
 2. calculate V_{λ^*} ; $\lambda \leftarrow \lambda^*$, $V_\lambda \leftarrow V_{\lambda^*}$
- until the segmentation no longer changes

$L_\lambda(O, V_\lambda)$ increases at every iteration of algorithm 2:

- At step 1, $L_\lambda(O, V_\lambda) \leq L_{\lambda^*}(O, V_\lambda)$: training the networks reduces their errors along V_λ ; the maximum likelihood estimates of the variances and transition probabilities are calculated with formulas 3 and 4.
- At step 2, $L_{\lambda^*}(O, V_\lambda) \leq L_{\lambda^*}(O, V_{\lambda^*})$ because V_{λ^*} is the most likely path for model λ .

We have: $L_\lambda(O, V_\lambda) \leq L_{\lambda^*}(O, V_\lambda) \leq L_{\lambda^*}(O, V_{\lambda^*})$.

3 Preliminary results

To understand what our induction algorithm can do, we first applied it to artificial data. A time series with 206 data points was generated by a process switching between two function patterns, f and g . For both patterns the self-transition probability is .9. Both functions are linear, f depends only on $x(t)$, g depends on both $x(t)$ and $x(t-1)$:

$$\begin{aligned} x(t+1) &= f(\cdot) = -.95x(t) + .5 \\ x(t+1) &= g(\cdot) = 1.05x(t) - .95x(t-1) - .5 \end{aligned}$$

Normally distributed noise with variance .01 was added to the output. The networks are trained on the $\langle x(t), x(t+1) \rangle$ pairs. Because the domains and codomains of the two functions overlap, the task of the induction algorithm is not trivial. Nevertheless, the induced segmentation, shown in figure 3, identifies almost perfectly the two patterns. The first induced network state gets all the 125 points generated by f plus 2 points generated by g , and the other three induced states get the rest of g 's points. All networks have two hidden units with cosine activation function and one state unit with identity activation function. Because f can be immediately computed from $x(t)$, a network approximating f was easily found by the algorithm. For g a network needs to save information about the previous value $x(t-1)$. Although each network in this experiment can approximate g well by storing $x(t-1)$ in its state unit, in several runs the algorithm found a three-network solution.

The algorithm was also applied to time series collected during experiences involving a robot approaching or passing an object. There were 14 kinds of experiences – “pass right a red object”, “pass right a red object, then push a blue object”, totaling 1082 time steps. Two thirds of the experiences were selected for the training set, and the remaining ones formed the test set. A model was induced from the training set for the task of predicting the next “vis-A-x” and “vis-A-y” sensor values from the current “trans-vel”, “vis-A-x” and “vis-A-y” values; “vis-A-x” and “vis-A-y” are the coordinates of a red object in the robot's visual field, and “trans-vel” is the translational velocity. Small networks with four hidden and two state cosine units were used. The segmentations of four “pass right a red object” experiences, two in the training set (top plots) and two in the test set

