

Hidden Markov Models

1 Introduction

[TBW]

2 Basic Hidden Markov Models

Hidden Markov Models (HMM) are created to model a sequence of events, based on hidden observations. In the demonstrated model in Fig. 1, let's assume that $\mathbf{X} = (X_1, \dots, X_n)$ is the set of observations, and $\mathbf{Z} = (Z_1, \dots, Z_n)$ is the set of latent states. Each state X_i could take on m_x possible value, and each Z_i could take on m_z possible values. The matrix $A_{m_z \times m_z}$ characterizes the transition probabilities for Z variables. Similarly the matrix $B_{m_z \times m_x}$ denote the probability of observing X_t from Z_t . So we are assuming the Markovian transitions from a hidden-state to a next hidden-states.

In this sequential model, we are in interested in the following probabilities:

- Inference: Given the sequence (\mathbf{x}, \mathbf{Z}) , what is the probability of $p(\mathbf{x}|\mathbf{Z})$.
- Learning: Given the sequence \mathbf{x} what is the sequence \mathbf{Z} that maximizes $p(\mathbf{Z}|\mathbf{x})$.

Now the question is given a sequence of training data, how can we find the most likely model which has generated the given sequence? We should assume an initial distribution for Z_1 by $w_{m_z \times 1}$. We can write the log-likelihood by

$$\mathcal{L} = \log w(Z_1) + \sum_{t=1}^{n-1} \log A(Z_t, Z_{t+1}) + \sum_{t=1}^n \log B(Z_t, X_t)$$

Figure 1: Structure of a Hidden Markov Model.

We shall train the model using EM algorithm, which also known as Baum-Welch. We define

$$\gamma_t(i, j) = p_{\theta_0}(Z_t = i, Z_{t+1} = j | \mathbf{x}), \quad \gamma_t(i) = p_{\theta_0}(Z_t = i | \mathbf{x}) = \sum_j \gamma_t(i, j)$$

E-Step: In terms of the above definitions,

$$\mathbb{E}_{\mathbf{Z}|\mathbf{x}, \theta} \log p(\mathbf{Z}, \mathbf{x} | \theta) = \sum_{i=1}^{m_z} \gamma_1(i) \log w(i) + \sum_{i,j=1}^{m_z} \sum_{t=1}^{n-1} \gamma_t(i, j) \log A(i, j) + \sum_{i=1}^{m_z} \sum_{t=1}^n \gamma_t(i) \log B(i, x_t)$$

M-Step: and using derivatives of the the above term we have the following estimated values:

$$\begin{aligned} w^*(i) &= \gamma_1(i), \quad i = 1, \dots, m_z; \\ A^*(i, j) &= \frac{\sum_{t=1}^{n-1} \gamma_t(i, j)}{\sum_{j'} \sum_{t=1}^{n-1} \gamma_t(i, j')}, \quad i, j = 1, \dots, m_z; \\ B^*(i, l) &= \frac{\sum_{t=1; x_t=l} \gamma_t(i)}{\sum_t \gamma_t(i)}, \quad i = 1, \dots, m_z, \quad l = 1, \dots, m_x; \end{aligned}$$

Forward-Backward probabilities: Define forward $\alpha_{t+1}(i)$ and backward $\beta_{t-1}(i)$ probabilities as following:

$$\alpha_{t+1}(i) = p_{\theta}(x_1, \dots, x_{t+1}, Z_{t+1} = i), \quad \beta_{t-1}(i) = p_{\theta}(x_t, \dots, x_n | Z_{t-1} = i)$$

We can show that the forward probability $\alpha_{t+1}(i)$ could be expressed based of the previous ones:

$$\alpha_{t+1}(i) = p_{\theta}(x_1, \dots, x_{t+1}, Z_{t+1} = i) \quad (1a)$$

$$= \sum_j p_{\theta}(x_1, \dots, x_{t+1}, Z_t = j, Z_{t+1} = i) \quad (1b)$$

$$= \sum_j \alpha_t(j) A(j, i) B(i, x_{t+1}) \quad (1c)$$

Similarly the backward probability could be written based on the future ones:

$$\beta_{t-1}(i) = p_{\theta}(x_t, \dots, x_n | Z_{t-1} = i) \quad (2a)$$

$$= \sum_j p_{\theta}(x_t, \dots, x_n, Z_t = j, Z_{t-1} = i) \quad (2b)$$

$$= \sum_j A(i, j) B(j, x_t) \beta_t(j) \quad (2c)$$

Now we can define the joint probabilities in terms of $\alpha(\cdot)$ and $\beta(\cdot)$.

$$p(Z_t = i, x_1, \dots, x_n) = p(x_1, \dots, x_n | Z_t = i) \cdot p(Z_t = i) \quad (3a)$$

$$= p(x_1, \dots, x_t, Z_t = i) \cdot p(x_{t+1}, \dots, x_n | Z_t = i) \quad (3b)$$

$$= \alpha_t(i) \cdot \beta_t(i) \quad (3c)$$

And also we can express the whole likelihood in terms of forward probabilities:

$$p(\mathbf{x}|\theta) = \sum_i p_\theta(x_1, \dots, x_n, Z_n = i) = \sum_i \alpha_n(i)$$

The computation of all probabilities in terms of forward/backward probabilities makes everything so simple. We just need to calculate them in two swipes. It turns out that forward-backward algorithms is a special case of *belief-propagation* algorithm in Graphical Models.

Now the question is given a model, how can we make inference about the hidden states. The inference about the hidden variables depend on the definition of the optimality for the inference. One criterion is to choose the hidden values that are individually most likely, that is:

$$\hat{Z}_t^* = \arg \max_i p(Z_t = i | \mathbf{x}, \theta) = \arg \max_i \gamma_t(i)$$

But the problem is that considering each observation individually is not reasonable, and one needs to choose the most-likely sequence:

$$\hat{\mathbf{Z}}^* = \arg \max_{i_1, \dots, i_n} p(Z_1 = i_1, \dots, Z_n = i_n | \mathbf{x}, \theta)$$

This solution could be found using Viterbi algorithm. Assume that we have a trained HMM and we want to find the most probable latent path, given some observations. The highest probability of the sequence of hidden states given first t observations is:

$$\delta_t(i) = \max_{j_1, \dots, j_{t-1}} p(Z_1 = j_1, \dots, Z_{t-1} = j_{t-1}, Z_t = i, \mathbf{x}_{1:t} | \theta)$$

Rewriting the optimality criterion in a recursive form :

$$\delta_{t+1}(i) = \left[\max_j \delta_t(j) A(j, i) \right] \cdot B(i, x_{t+1})$$

Similar implementation to forward/backward algorithm, will give us the optimal latent sequence. This turns out that Viterbi is a special case of *max-product* algorithm in graphical models.