PGM 2004/05 Tirgul 2

Hidden Markov Models

based on Gal Elidan notes

Introduction

Hidden Markov Models (HMM) are one of the most common form of probabilistic graphical models, although they were developed long before the notion of general models existed (1913). They are used to model time-invariant and limited horizon models that have both an underlying mechanism (hidden states) and an observable consequence. They have been extremely successful in language modeling and speech recognition systems and are still the most widely used technique in these domains.

Markov Models

A Markov process or model assumes that we can predict the future based just on the present (or on a limited horizon into the past):

Let $\{X_1, ..., X_T\}$ be a sequence of random variables taking values $\{1, ..., N\}$ then the Markov properties are:

Limited Horizon:

$$P(X_{t+1} = i | X_1, ..., X_t) = P(X_{t+1} = i | X_t)$$

Time invariant (stationary):

$$= P(X_2 = i | X_1)$$

Describing a Markov Chain

A Markov chain can be described by the transition matrix A and initial state probabilities Q:

$$A_{ij} = P(X_{t+1} = j | X_t = i)$$
 $q_i = P(X_1 = i)$

or alternatively:



and we calculate:

 $P(X_1,...,X_T) = P(X_1)P(X_2 \mid X_1)...P(X_T \mid X_{T-1}) = q_{X_1}\prod_{t=1}^{T-1}A(X_t,X_{t+1})$

Hidden Markov Models

In a Hidden Markov Model (HMM) we do not observe the sequence that the model passed through but only some probabilistic function of it. Thus, it is a Markov model with the addition of *emission probabilities*:

$$B_{ik} = P(Y_t = k \mid X_t = i)$$

Why use HMMs?

• A lot of real-life processes are composed of underlying events generating surface phenomena. Tagging parts of speech is a common example.

• We can usually think of processes as having a limited horizon (we can easily extend to the case of a constant horizon larger than 1)

- We have efficient training algorithm using EM
- Once the model is set, we can easily run it:

t=1, start in state i with probability q_i

forever : move from state i to j with probability a_{ii}

emit y_t =k with probability b_{ik}

t=t+1

The fundamental questions

Likelihood: Given a model $\lambda = (A,B,q)$, how do we efficiently compute the likelihood of an observation $P(Y \mid \lambda)$?

Decoding: Given the observation sequence Y and a model μ , what state sequence explains it best (MPE)?

This is, for example, the tagging process of an observed sentence.

Learning: Given an observation sequence Y, and a generic model, how do we estimate the parameters that define the best model to describe the data?

Computing the Likelihood

Using P(Y,X)=P(Y|X)P(X) we get:

$$P(Y) = \sum_{X} P(Y,X) = \sum_{X} P(Y|X)P(X)$$

$$= \sum_{x_1 \cdots x_T} q_{x_1} \prod_{t=2}^T A_{x_{t-1}x_t} B_{x_t y_t}$$

But, we have $O(TN^{T})$ multiplications!

The trellis (lattice) algorithm

To compute likelihood: Need to enumerate over all paths in the lattice (all possible instantiations of $X_1...X_T$).

But... some starting subpath (blue) is common to many continuing paths (blue+red)



Idea: using dynamic programming, calculate a path in terms of shorter sub-paths

The trellis (lattice) algorithm

We build a matrix of the probability of being at time t at state i - $\alpha_t(i)=P(x_t=i,y_1y_2...y_t)$ is a function of the previous column (forward procedure):



The trellis (lattice) algorithm (cont.)

We can similarly define a backwards procedure for filling the matrix $\beta_{t}(i) = P(\gamma_{t+1} \cdots \gamma_{T} \mid X_{t} = i)$



The trellis (lattice) algorithm (cont.) And we can easily combine:

$$P(Y, X_{t} = i) = P(y_{1} \cdots y_{T}, X_{t} = i)$$

= $P(y_{1} \cdots y_{t}, X_{t} = i, y_{t+1} \cdots y_{T})$
= $P(y_{1} \cdots y_{t}, X_{t} = i)P(y_{t+1} \cdots y_{T} | y_{1} \cdots y_{t}, X_{t} = i)$
= $P(y_{1} \cdots y_{t}, X_{t} = i)P(y_{t+1} \cdots y_{T} | X_{t} = i)$
= $\alpha_{t}(i)\beta_{t}(i)$

And then

$$\mathsf{P}(\mathsf{Y}) = \sum_{i} \mathsf{P}(\mathsf{Y}, \mathsf{X}_{\dagger} = i) = \sum_{i} \alpha_{\dagger}(i)\beta_{\dagger}(i)$$

Finding the best state sequence

We would like to the most likely path (and not just the most likely state at each time slice)

$$\underset{x}{\operatorname{argmax}} P(X \mid Y) = \underset{x}{\operatorname{argmax}} P(X,Y)$$

The Viterbi algorithm is an efficient trellis method for finding the MPE:

$$\delta_1(i) = q_i B_{iy_1}$$

$$\delta_{t+1}(i) = B_{iy_{t+1}} \max_j \delta_t(j) A_{ji} \qquad \gamma_{t+1}(i) = \arg\max_j \delta_t(j) A_{ji}$$

and we to reconstruct the path:

$$P(\hat{X}) = \max_{i} \delta_{T}(i) \quad \hat{X}_{T} = \arg_{i} \max \delta_{T}(i)$$
$$\hat{X}_{t-1} = \gamma_{t}(\hat{X}_{t})$$

The Casino HMM

A casino switches from a fair die (state F) to a loaded one (state U) with probability 0.05 and the other way around with probability 0.1. The loaded die has 0.5 probability of showing a 6. The casino, honestly, reads off the number that was rolled.

$$\boldsymbol{\mathcal{A}} = \begin{pmatrix} 0.95 & 0.05 \\ 0.1 & 0.9 \end{pmatrix} \qquad \boldsymbol{\mathcal{B}} = \begin{pmatrix} 1/4 & 1/4 & 1/4 & 1/4 & 1/4 & 1/4 \\ 1/6 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 \\ 1/10 & 1/10 & 1/10 & 1/10 & 1/2 \end{pmatrix}$$

 $P(3151166661 | FFFFUUUUU) = \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{2} \cdot \frac{1}{2$

The Casino HMM (cont.)

What is the likelyhood of 3151166661?

- Y= 3151166661
- $\alpha_1(1)=0.5*1/6=1/12,$

 $\alpha_1(2)=0.5*0.1=0.05$

- $\alpha_2(1)=1/6^*(0.95^*1/12+0.1^*0.05)\cong 0.014$
- $\alpha_2(1)=0.1^*(0.05^*1/12+0.9^*0.05)\cong 0.0049$
- $\alpha_3(1) \cong 0.0023, \ \alpha_3(2) \cong 0.0005$
- $\alpha_4(1) \cong 0.0004$, $\alpha_4(1) \cong 0.0001$
- $\alpha_5(1) \cong 0.0001$, $\alpha_5(1) < 0.0001$
- ... all smaller then 0.0001!

The Casino HMM (cont.)

What explains 3151166661 best?

```
Y= 3151166661
```

 $\delta_1(1)=0.5*1/6=1/12, \delta_1(2)=0.5*0.1=0.05$

 $\delta_2(1){=}1/6^*max(0.95^*1/12,0.1^*0.05)\cong 0.0132$

 $\delta_2(1)=1/10^* \max(0.05^*1/12, 0.9^*0.05) \cong 0.0045$

 $\delta_3(1) \cong 0.0021, \ \delta_3(2) \cong 0.0004$

 $\delta_4(1) \cong 0.0003, \, \delta_4(1) < 0.0001$

 $\delta_5(1) \cong 0.0001, \, \delta_5(1) < 0.001 \dots$

. . .

The Casino HMM (cont.)

An example of reconstruction using Viterbi (Durbin):

Rolls 3151162464466442453113216311641521336 Die Viterbi 2514454363165662656666665116645313265 Rolls 00000001111111111111111111111110000000 Die Viterbi 00000000001111111111111111111110000000 1245636664631636663162326455236266666 Rolls Die Viterbi

Learning

If we were given both X and Y, we could choose

$$\underset{A,B}{\operatorname{AB}} \operatorname{AB} P(Y_{training}, X_{training} \mid A, B, Q)$$

Using the Maximum Likelihood principal, we simply assign the parameter for each relative frequency

What do we do when we have only Y?

$$\underset{A,B}{\operatorname{argmax}} P(Y_{training} \mid A, B, Q)$$

ML here does not have a closed form formula!

EM (Baum Welch)

Idea: Using current guess to complete data and reestimate



Thm: Likelihood of observables never decreases!!! (to be proved later in the course)

Problems: Gets stuck at sub-optimal solutions

Parameter Estimation

We define the expected number of transitions from state i to j at time t:

$$p_t(i,j) = P(X_t = i, X_{t+1} = j | Y) = \frac{\alpha_t(i)A_{ij}B_{jy_{t+1}}\beta_{t+1}(j)}{\sum_m \sum_n \alpha_t(m)A_{mn}B_{ny_{t+1}}\beta_{t+1}(n)}$$

the expected # of transition from i to j in Y is then

$$\sum_{t=1}^{T} p_t(i,j)$$

Parameter Esimation (cont.)

We use EM re-estimation formulas using the expected counts we already have:

$$\overline{\mathbf{q}_{i}} = \exp \operatorname{ected} \operatorname{frequency} \operatorname{in} \operatorname{state} \operatorname{i} \operatorname{at} \operatorname{time} 1 = \sum_{j} P_{1}(i,j)$$

$$\overline{\mathbf{a}_{ij}} = \frac{\exp \operatorname{ected} \# \operatorname{of} \operatorname{transitions} \operatorname{from} \operatorname{state} \operatorname{i} \operatorname{to} j}{\exp \operatorname{ected} \# \operatorname{of} \operatorname{transitions} \operatorname{from} \operatorname{state} i} = \frac{\sum_{j} p_{t}(i,j)}{\sum_{j} p_{t}(i,j)}$$

$$\overline{\mathbf{b}_{j}} = \exp \operatorname{ected} \# \operatorname{of} \operatorname{emissions} \operatorname{of} k \operatorname{from} \operatorname{state} i = \sum_{j} \sum_{\substack{j \ t: y_{t} = k}} p_{t}(i,j)$$

) _p_t(i,j)

exp ected # of emissions from state i

Dik