A gentle introduction to Hidden Markov Models

> Mark Johnson Brown University

November 2009

#### What is sequence labeling?

- Markov models
- Hidden Markov models
- Finding the most likely state sequence
- Estimating HMMs
- Conclusion

# Sequence labeling

- Input: a sequence  $\boldsymbol{x} = (x_1, \dots, x_n)$  (usually n may vary)
- Output: a sequence  $y = (x_1, \ldots, y_n)$ , where  $y_i$  is a label for  $x_i$

Part of speech (POS) tagging:

- $\boldsymbol{y}$ : DT JJ NN VBD NNP.
- $\boldsymbol{x}$ : the big cat bit Sam .

Noun-phrase chunking:

$oldsymbol{y}$ :	[NP	NP	NP]	-	[NP]	
$oldsymbol{x}$ :	the	big	$\operatorname{cat}$	bit	$\operatorname{Sam}$	

Named entity detection:

# Sequence labeling with probabilistic models

- General idea: develop a probabilistic model  $\operatorname{P}(oldsymbol{y},oldsymbol{x})$ 
  - use this to compute  $\mathrm{P}(m{y}|m{x})$  and  $\hat{m{y}}(m{x}) = \mathrm{argmax}_{m{y}} \mathrm{P}(m{y}|m{x})$
- A Hidden Markov Model (HMM) factors:

 $P(\boldsymbol{y}, \boldsymbol{x}) = P(\boldsymbol{y}) P(\boldsymbol{x}|\boldsymbol{y})$ 

- $\mathrm{P}(oldsymbol{y})$  is the probability of sequence  $oldsymbol{y}$ 
  - In an HMM, P(y) is a *Markov chain* 
    - in applications,  $oldsymbol{y}$  is usually "hidden"
  - P(x|y) is the *emission model*.

#### What is sequence labeling?

- Markov models
- Hidden Markov models
- Finding the most likely state sequence
- Estimating HMMs
- Conclusion

Markov models of sequences

• Let  $\boldsymbol{y} = (y_1, \dots, y_n)$ , where  $y_i \in \mathcal{Y}$  and n is fixed  $P(\boldsymbol{y}) = P(y_1) P(y_2|y_1) P(y_3|y_1, y_2) \dots P(y_n|y_1, \dots, y_{n-1})$ 

= 
$$P(y_1) \prod_{i=2}^{n} P(y_i|y_{1:i-1})$$

where  $y_{i:j}$  is the subsequence  $(y_i, \ldots, y_j)$  of y

• In a (first-order) Markov chain:

$$P(y_i|y_{1:i-1}) = P(y_i|y_{i-1}) \text{ for } i > 1$$

In a *homogenous* Markov chain, P(y<sub>i</sub>|y<sub>i-1</sub>) does not depend on *i*, so probability of a sequence is completely determined by:

$$P(y_1) = \iota(y_1) \quad (start \ probabilities)$$

$$P(y_i|y_{i-1}) = \sigma(y_{i-1}, y_i) \quad (transition \ probabilities), \text{ so:}$$

$$P(y) = \iota(y_1) \prod_{i=2}^n \sigma(y_{i-1}, y_i)$$

# Markov models of varying-length sequences

- To define a model of sequences y = (y<sub>1</sub>,..., y<sub>n</sub>) of varying lengths n, need to define P(n)
- An easy way to do this:
  - add a new *end marker* symbol ' $\triangleleft$ ' to  ${\mathcal Y}$
  - pad all sequences with this end-marker

$$\boldsymbol{y} = y_1, y_2, \ldots, y_n, \triangleleft, \triangleleft, \ldots$$

• set  $\sigma(\triangleleft, \triangleleft) = 1$ , i.e.,  $\triangleleft$  is a *sink state* 

- ▶ so  $\sigma(y, \triangleleft)$  is probability of sequence ending after state y
- We can also incorporate the start probabilities *ι* into *σ* by introducing a new *begin marker* y<sub>0</sub> = ▷

$$\boldsymbol{y} = \triangleright, y_1, y_2, \dots, y_n, \triangleleft, \triangleleft, \dots, \text{ so:}$$
  

$$P(\boldsymbol{y}) = \prod_{i=1}^{\infty} \sigma(y_{i-1}, y_i) \text{ (only compute to } i = n+1)$$

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

## Hidden Markov models

- An HMM is specified by:
  - A state-to-state transition matrix σ, where σ(y, y') is the probability of moving from state y to state y'
  - A state-to-observation emission matrix τ, where τ(y, x) is the probability of emitting x in state y

$$P(\boldsymbol{y}, \boldsymbol{x}) = P(\boldsymbol{y}) P(\boldsymbol{x}|\boldsymbol{y})$$

$$P(\boldsymbol{y}) = \prod_{i=1}^{n+1} P(y_i|y_{i-1}) = \prod_{i=1}^{n+1} \sigma(y_{i-1}, y_i)$$

$$P(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i=1}^{n} P(x_i|y_i) = \prod_{i=1}^{n} \tau(y_i, x_i)$$

• I.e., each  $x_i$  is generated conditioned on  $y_i$  with probability  $\tau(y_i, x_i)$ 

#### HMMs as stochastic automata



Bayes net representation of an HMM

Trellis representation of HMM



 $X_1 = flour$   $X_2 = pan$ 

- The trellis representation represents the possible state sequences  $m{y}$  for a single observation sequence  $m{x}$
- A trellis is a DAG with *nodes* (i, y), where  $i \in 0, ..., n+1$  and  $y \in \mathcal{Y}$  and *edges* from (i 1, y') to (i, y)
- It plays a crucial role in *dynamic programming algorithms* for HMMs

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

#### Two important inference problems

- Given an HMM  $\sigma, au$  and an observation sequence x:
  - what is the most likely state sequence  $\hat{y} = \operatorname{argmax}_{y} P(x, y)$ ?
  - what is the probability of a sequence  $P(x) = \sum_{y} P(x, y)$ ?
- If the sequences are short enough, can be solved by exhaustively enumerating  $\boldsymbol{y}$
- Given a sequence of length n and an HMM with m states, can be solved using dynamic programming in  $O(n m^2)$  time

Key observation:  $P(\boldsymbol{x}, \boldsymbol{y})$  factorizes!

$$P(\boldsymbol{x}, \boldsymbol{y}) = \left(\prod_{i=1}^{n} \sigma(y_{i-1}, y_i) \tau(y_i, x_i)\right) \sigma(y_n, \triangleleft), \text{ so:} \\ -\log P(\boldsymbol{x}, \boldsymbol{y}) = \left(\sum_{i=1}^{n} -\log(\sigma(y_{i-1}, y_i) \tau(y_i, x_i))\right) - \log \sigma(y_n, \triangleleft)$$

• Given observation sequence x, define a trellis with *edge weights* 

$$w((i-1,y'), (i,y)) = -\log(\sigma(y_{i-1},y_i)\tau(y_i,x_i)) w((n,y), (n+1, \triangleleft)) = -\log(\sigma(y, \triangleleft))$$

 Then the - log probability of any state sequence y is the sum of the weights of the corresponding path in the trellis

# Finding the most likely state sequence $\hat{m{y}}$

- Goal: find  $\hat{m{y}} = \mathrm{argmax}_{m{y}} \operatorname{P}(m{x},m{y})$
- Algorithm:
  - Construct a trellis with edge weights

$$w((i-1,y'), (i,y)) = -\log(\sigma(y_{i-1},y_i)\tau(y_i,x_i))$$
  
$$w((n,y), (n+1, d)) = -\log(\sigma(y, d))$$

Run a shortest path algorithm on this trellis, which returns a state sequence with lowest – log probability

# The Viterbi algorithm for finding $\hat{y}$

Idea: given x = (x<sub>1</sub>,..., x<sub>n</sub>), let μ(k, u) be the maximum probability of any state sequence y<sub>1</sub>,..., y<sub>k</sub> ending in y<sub>k</sub> = u ∈ 𝔅.

$$u(k, u) = \max_{\boldsymbol{y}_{1:k}: y_k=u} P(\boldsymbol{x}_{1:k}, \boldsymbol{y}_{1:k})$$

• Base case: 
$$\mu(0, \triangleright) = 1$$

Recursive case:

$$\mu(k, u) = \max_{u'} \mu(k - 1, u') \sigma(u', u) \tau(u, x_k)$$

Final case:  $\mu(n+1, \triangleleft) = \max_{u'} \mu(n, u') \sigma(u', \triangleleft)$ 

• To find  $\hat{y}$ , for each trellis node (k, u) keep a *back-pointer* to *most likely predecessor* 

$$\rho(k, u) = \operatorname{argmax}_{u' \in \mathcal{Y}} \mu(k - 1, u') \sigma(u', u) \tau(u, x_k)$$

# Why does the Viterbi algorithm work?

$$\mu(k, u) = \max_{y_{1:k}: y_{k}=u} P(x_{1:k}, y_{1:k})$$

$$= \max_{y_{1:k}: y_{k}=u} \prod_{i=1}^{k} \sigma(y_{i-1}, y_{i}) \tau(y_{i}, x_{i})$$

$$= \max_{y_{1:k-1}} \left( \prod_{i=1}^{k-1} \sigma(y_{i-1}, y_{i}) \tau(y_{i}, x_{i}) \right) \sigma(y_{k-1}, u) \tau(u, x_{k})$$

$$= \max_{u'} \left( \max_{y_{1:k-1}: y_{k-1}=u'} \prod_{i=1}^{k-1} \sigma(y_{i-1}, y_{i}) \tau(y_{i}, x_{i}) \right) \sigma(u', u) \tau(u, x_{k})$$

$$= \max_{u'} \mu(k-1, u') \sigma(u', u) \tau(u, x_{k})$$

# The sum-product algorithm for computing P(x)

- Goal: compute  $P(\boldsymbol{x}) = \sum_{\boldsymbol{y} \in \mathcal{Y}^n} P(\boldsymbol{x}, \boldsymbol{y})$
- Idea: given x = (x<sub>1</sub>,...,x<sub>n</sub>), let α(k, u) be the sum of the probabilities of x<sub>1</sub>,..., x<sub>k</sub> and all state sequences y<sub>1</sub>,..., y<sub>k</sub> ending in y<sub>k</sub> = u.

$$\alpha(k, u) = \sum_{y_{1:k}: y_k=u} P(x_{1:k}, y_{1:k})$$

• Recursive case: for  $k = 1, \ldots, n$ :

$$\alpha(k, u) = \sum_{u'} \alpha(k-1, u') \sigma(u', u) \tau(u, x_i)$$

• Final case:  $\alpha(n+1, \triangleleft) = \sum_{u'} \alpha(n, u') \sigma(u', \triangleleft)$ 

• The lpha are called *forward probabilities* 

What is sequence labeling?

Markov models

Hidden Markov models

Finding the most likely state sequence

Estimating HMMs

Conclusion

## Estimating HMMs from visible data

- Training data:  $D = ((x_1, y_1), \dots, (x_n, y_n))$ , where the  $x_i$  and  $y_i$  are *sequences* 
  - ▶ for conceptual simplicity, assume that the data comes as two long sequences D = (x, y)
- Then the MLE is:

$$\begin{aligned} \hat{\sigma}(y',y) &= \frac{c(y',y)}{\sum_{y\in\mathcal{Y}}c(y',y)} \\ \hat{\tau}(y,x) &= \frac{c(y,x)}{\sum_{x\in\mathcal{X}}c(y,x)}, \text{ where:} \\ c(y',y) &= \text{ the number of times state } y' \text{ precedes state } y \text{ in } y \\ c(y,x) &= \text{ the number of times state } y \text{ emits } x \text{ in } (x,y) \end{aligned}$$

# Estimating HMMs from *hidden data* using EM

• EM iteration at time t:

$$\sigma(y',y)^{(t)} = \frac{\mathrm{E}[c(y',y)]}{\sum_{y\in\mathcal{Y}}\mathrm{E}[c(y',y)]}$$
  

$$\tau(y,x)^{(t)} = \frac{\mathrm{E}[c(y,x)]}{\sum_{x\in\mathcal{X}}\mathrm{E}[c(y,x)]}, \text{ where:}$$
  

$$\mathrm{E}[c(y',y)] = \sum_{y\in\mathcal{Y}^n} c(y',y) \mathrm{P}(y|x,\sigma^{(t-1)},\tau^{(t-1)})$$
  

$$\mathrm{E}[c(y,x)] = \sum_{y\in\mathcal{Y}^n} c(y,x) \mathrm{P}(y|x,\sigma^{(t-1)},\tau^{(t-1)})$$

- Key computational problem: *computing these expectations efficiently*
- Dynamic programming to the rescue!

## Expectations require marginal probabilities

• Suppose we could efficiently compute the *marginal probabilities*:

$$P(Y_i = y'|x) = \frac{\sum_{y \in \mathcal{Y}^n : y_i = y'} P(y, x)}{P(x)}$$

• Then we could efficiently compute:

$$\mathbf{E}[c(y',x')] = \sum_{i:x_i=x'} \mathbf{P}(Y_i = y'|x)$$

Efficiently computing E[c(y', y)] requires computation of a different marginal probability

#### Backward probabilities

Idea: given x = (x<sub>1</sub>,...,x<sub>n</sub>), let β(k, v) be the sum of the probabilities of x<sub>k+1</sub>,..., x<sub>n</sub> and all state sequences y<sub>k+1</sub>,..., y<sub>n</sub> conditioned on y<sub>k</sub> = v.

$$eta(k, \mathbf{v}) = \sum_{y_{k:n}: y_k = \mathbf{v}} \mathrm{P}(\boldsymbol{x}_{k+1:n}, \boldsymbol{y}_{k+1:n} | Y_k = \mathbf{v})$$

- ▶ Base case:  $\beta(n, v) = \sigma(v, \triangleleft)$  for each  $v \in \mathcal{Y}$
- Recursive case: for  $k = n 1, \ldots, 0$ :

$$\beta(k, \mathbf{v}) = \sum_{\mathbf{v}'} \sigma(\mathbf{v}, \mathbf{v}') \tau(\mathbf{v}', \mathbf{x}_{k+1}) \beta(k+1, \mathbf{v}')$$

• The  $\beta$  are called *backward probabilities* 

# Computing marginals from forward and backward probabilities

$$P(Y_{i} = u | \boldsymbol{x}) = \frac{1}{P(\boldsymbol{x})} \sum_{\boldsymbol{y} \in \mathcal{Y}^{n} : y_{i} = u} P(\boldsymbol{y}, \boldsymbol{x})$$

$$= \frac{1}{P(\boldsymbol{x})} \sum_{y_{1:k} : y_{i} = u} \sum_{y_{i:n} : y_{i} = u} P(\boldsymbol{x}_{1:i}, \boldsymbol{y}_{1:i}) P(\boldsymbol{x}_{i+1:n}, \boldsymbol{y}_{i:n})$$

$$= \frac{1}{P(\boldsymbol{x})} \left( \sum_{y_{1:k} : y_{i} = u} P(\boldsymbol{x}_{1:i}, \boldsymbol{y}_{1:i}) \right) \left( \sum_{y_{k:n} : y_{i} = u} P(\boldsymbol{x}_{i+1:n}, \boldsymbol{y}_{i:n}) \right)$$

$$= \frac{1}{P(\boldsymbol{x})} \alpha(i, u) \beta(i, u)$$

• So the marginal  $\mathrm{P}(Y_i|x)$  can be efficiently computed from lpha and eta

- What is sequence labeling?
- Markov models
- Hidden Markov models
- Finding the most likely state sequence
- Estimating HMMs
- Conclusion

# Conclusion

- Hidden Markov Models (HMMs) can be used to label a sequence of observations
- There are efficient dynamic programming algorithms to find the most likely label sequence and the marginal probability of a label
- The expectation-maximization algorithm can be used to learn an HMM from unlabeled data
- The expected values required for EM can be efficiently computed