

# Modeling Sequence Data: HMMs and Viterbi

CS4780/5780 – Machine Learning  
Fall 2014

Tobias Schnabel and Igor Labutov  
Cornell University

Reading:  
Manning/Schütze, Sections 9.1-9.3 (except 9.3.1)  
Leeds Online HMM Tutorial (except Forward and Forward/Backward Algorithm)  
([http://www.comp.leeds.ac.uk/roger/HiddenMarkovModels/html\\_dev/main.html](http://www.comp.leeds.ac.uk/roger/HiddenMarkovModels/html_dev/main.html))

## Outline

- Hidden Markov Models
  - Viterbi Algorithm
  - Estimation with fully observed training data
  - Applications: Part-of-speech tagging



## Hidden Markov Model

- States:  $y \in \{s_1, \dots, s_k\}$
- Outputs symbols:  $x \in \{o_1, \dots, o_m\}$

Parameter	
Starting probability	$P(Y_1 = y_1)$
Transition probability	$P(Y_i = y_i   Y_{i-1} = y_{i-1})$
Output/Emission probability	$P(X_i = x_i   Y_i = y_i)$

## Hidden Markov Model

- Every output/state sequence has a probability

$$\begin{aligned}
 P(x, y) &= P(x_1, \dots, x_l, y_1, \dots, y_l) \\
 &= P(y_1)P(x_1|y_1) \prod_{i=2}^l P(x_i|y_i)P(y_i|y_{i-1})
 \end{aligned}$$

- Different visualizations

## Estimating the Probabilities

- Fully observed data:

$$\begin{aligned}
 P(Y_i = a | Y_{i-1} = b) &= \frac{\text{\# of times state } a \text{ follows state } b}{\text{\# of times state } b \text{ occurs}} \\
 P(X_i = a | Y_i = b) &= \frac{\text{\# of times output } a \text{ is observed in state } b}{\text{\# of times state } b \text{ occurs}}
 \end{aligned}$$

- Smoothing the estimates:

- See Naïve Bayes for text classification

- Partially observed data ( $Y_i$  unknown):
  - Expectation-Maximization (EM)

## HMM Decoding: Viterbi Algorithm

- Question: What is the most likely state sequence given an output sequence

– Find  $y^* = \operatorname{argmax}_{y \in \{y_1, \dots, y_l\}} P(x_1, \dots, x_l, y_1, \dots, y_l)$

$$= \operatorname{argmax}_{y \in \{y_1, \dots, y_l\}} \left\{ P(y_1) P(x_1 | y_1) \prod_{i=2}^l P(x_i | y_i) P(y_i | y_{i-1}) \right\}$$

## Going on a trip

- Deal: 3 trips to cities 3 different countries:



## Going on a trip

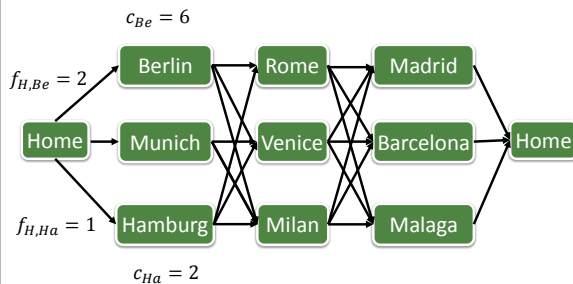
- Deal: 3 trips to cities 3 different countries:

Country	City options
Germany	Berlin/Munich/Hamburg
Italy	Rome/Venice/Milan
Spain	Madrid/Barcelona/Malaga

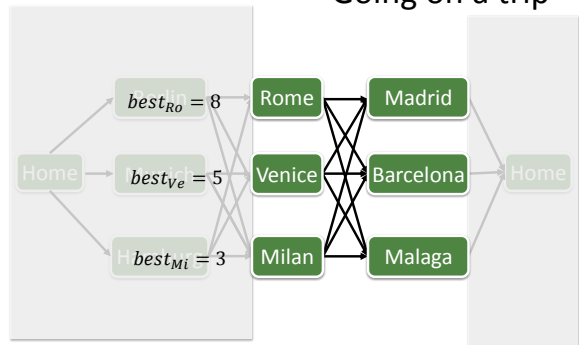
## Going on a trip

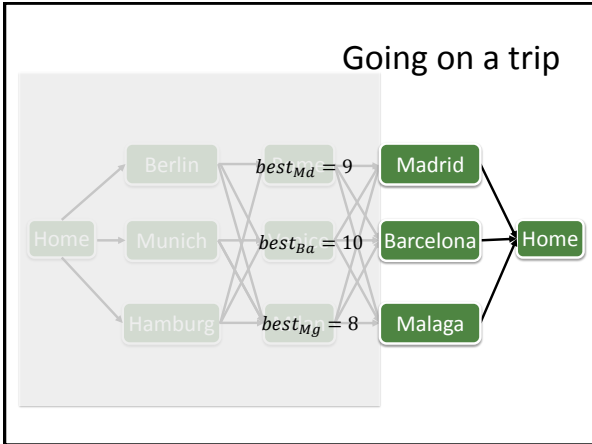
- Deal: 3 trips to cities 3 different countries:
  - Each city  $i$  has an attractiveness score  $c_i \in [0, 10]$
  - Each flight has a comfort score  $f_{i,j} \in [0, 10]$
- Find the best trip!

## Going on a trip



## Going on a trip





### HMM Decoding: Viterbi Algorithm

- Question: What is the most likely state sequence given an output sequence
  - Find  $y^* = \operatorname{argmax}_{y \in \{y_1, \dots, y_l\}} P(x_1, \dots, x_l, y_1, \dots, y_l)$
  - $$= \operatorname{argmax}_{y \in \{y_1, \dots, y_l\}} \left\{ P(y_1) P(x_1 | y_1) \prod_{i=2}^l P(x_i | y_i) P(y_i | y_{i-1}) \right\}$$
  - Viterbi algorithm has runtime linear in length of sequence

### Viterbi Example

$P(X_i   Y_i)$	A+	B	C
happy	0.6	0.3	0.1
grumpy	0.1	0.4	0.5

$P(Y_i)$		$P(Y_i   Y_{i-1})$	happy	grumpy
happy	0.7		0.8	0.2
grumpy	0.3		0.3	0.7

- What the most likely mood sequence for  $x = (C, A+, A+)$ ?

### HMM's for POS Tagging

- Design HMM structure (vanilla)
  - States: one state per POS tag
  - Transitions: fully connected
  - Emissions: all words observed in training corpus
- Estimate probabilities
  - Use corpus, e.g. Treebank
  - Smoothing
  - Unseen words?
- Tagging new sentences
  - Use Viterbi to find most likely tag sequence

### Experimental Results

Tagger	Accuracy	Training time	Prediction time
HMM	96.80%	20 sec	18.000 words/s
TBL Rules	96.47%	9 days	750 words/s

- Experiment setup
  - WSJ Corpus
  - Trigram HMM model
  - from [Pla and Molina, 2001]

### Discriminative vs. Generative

- Bayes Rule: 
$$h_{\text{bayes}}(x) = \operatorname{argmax}_{y \in Y} [P(Y = y | X = x)]$$

$$= \operatorname{argmax}_{y \in Y} [P(X = x | Y = y) P(Y = y)]$$
- Generative:
  - Model  $P(X = x | Y = y)$  and  $P(Y = y)$
- Discriminative:
  - Find  $h$  in  $H$  that best approximates the classifications made by
 
$$h_{\text{bayes}}(x) = \operatorname{argmax}_{y \in Y} [P(Y = y | X = x)]$$
- Question: Can we train HMM's discriminately?