

CSCI 5832
Natural Language Processing

Jim Martin
Lecture 9

2/28/08 1

Today 2/19

- Review HMMs for POS tagging
- Entropy intuition
- Statistical Sequence classifiers
 - ♦ HMMs
 - ♦ MaxEnt
 - ♦ MEMMs

2/28/08 2

Statistical Sequence Classification

- Given an input sequence, assign a label (or tag) to each element of the tape
 - ♦ Or... Given an input tape, write a tag out to an output tape for each cell on the input tape
- Can be viewed as a classification task if we view
 - ♦ The individual cells on the input tape as things to be classified
 - ♦ The tags written on the output tape as the class labels

2/28/08 3

POS Tagging as Sequence Classification

- We are given a sentence (an “observation” or “sequence of observations”)
 - ♦ *Secretariat is expected to race tomorrow*
- What is the best sequence of tags which corresponds to this sequence of observations?
- Probabilistic view:
 - ♦ Consider all possible sequences of tags
 - ♦ Out of this universe of sequences, choose the tag sequence which is most probable given the observation sequence of n words $w_1 \dots w_n$.

2/28/08

4

Statistical Sequence Classification

- We want, out of all sequences of n tags $t_1 \dots t_n$ the single tag sequence such that $P(t_1 \dots t_n | w_1 \dots w_n)$ is highest.

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- Hat ^ means “our estimate of the best one”
- $\operatorname{Argmax}_x f(x)$ means “the x such that f(x) is maximized”

2/28/08

5

Road to HMMs

- This equation is guaranteed to give us the best tag sequence

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- But how to make it operational? How to compute this value?
- Intuition of Bayesian classification:
 - ♦ Use Bayes rule to transform into a set of other probabilities that are easier to compute



2/28/08

6

Using Bayes Rule

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n) P(t_1^n)}{P(w_1^n)}$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n) P(t_1^n)$$

2/28/08

7

Likelihood and Prior

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \overbrace{P(w_1^n | t_1^n)}^{\text{likelihood}} \overbrace{P(t_1^n)}^{\text{prior}}$$

$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i)$$

$$P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$



$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

2/28/08

8

Transition Probabilities

- Tag transition probabilities $p(t_i | t_{i-1})$
 - ♦ Determiners likely to precede adjs and nouns
 - That/DT flight/NN
 - The/DT yellow/JJ hat/NN
 - So we expect $P(NN|DT)$ and $P(JJ|DT)$ to be high
 - ♦ Compute $P(NN|DT)$ by counting in a labeled corpus:

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

$$P(NN|DT) = \frac{C(DT, NN)}{C(DT)} = \frac{56,509}{116,454} = .49$$

2/28/08

9

Observation Probabilities

- Word likelihood probabilities $p(w_i|t_i)$
 - ♦ VBZ (3sg Pres verb) likely to be “is”
 - ♦ Compute $P(\text{is}|VBZ)$ by counting in a labeled corpus: $P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$

$$P(\text{is}|VBZ) = \frac{C(VBZ, \text{is})}{C(VBZ)} = \frac{10,073}{21,627} = .47$$

2/28/08

10

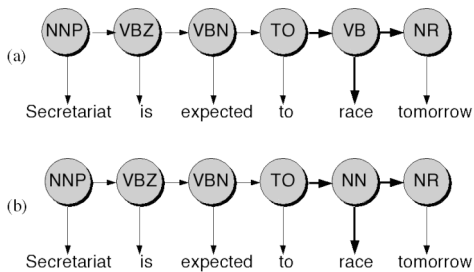
An Example: the verb “race”

- Secretariat/NNP is/VBZ expected/VBN to/TO **race**/VB tomorrow/NR
- People/NNS continue/VB to/TO inquire/VB the/DT reason/NN for/IN the/DT **race**/NN for/IN outer/JJ space/NN
- How do we pick the right tag?

2/28/08

11

Disambiguating “race”



2/28/08

12

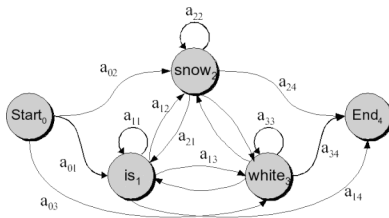
Example

- $P(\text{NN}|\text{TO}) = .00047$
- $P(\text{VB}|\text{TO}) = .83$
- $P(\text{race}|\text{NN}) = .00057$
- $P(\text{race}|\text{VB}) = .00012$
- $P(\text{NR}|\text{VB}) = .0027$
- $P(\text{NR}|\text{NN}) = .0012$
- $P(\text{VB}|\text{TO})P(\text{NR}|\text{VB})P(\text{race}|\text{VB}) = .00000027$
- $P(\text{NN}|\text{TO})P(\text{NR}|\text{NN})P(\text{race}|\text{NN}) = .0000000032$
- So we (correctly) choose the verb reading,

2/28/08

13

Markov chain for words



2/28/08

14

Markov chain = “First-order Observable Markov Model”

- A set of states
 - ♦ $Q = q_1, q_2, \dots, q_N$, the state at time t is q_t
- Transition probabilities:
 - ♦ a set of probabilities $A = a_{01}a_{02} \dots a_{n1} \dots a_{nn}$.
 - ♦ Each a_{ij} represents the probability of transitioning from state i to state j
 - ♦ The set of these is the transition probability matrix A
- Current state only depends on previous state

$$P(q_t = j | q_{t-1} = i) = a_{ij}$$

2/28/08

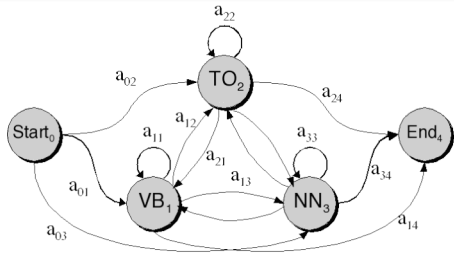
15

Hidden Markov Models

- States $Q = q_1, q_2, \dots, q_N$;
- Observations $O = o_1, o_2, \dots, o_N$;
 - Each observation is a symbol from a vocabulary $V = \{v_1, v_2, \dots, v_V\}$
- Transition probabilities
 - Transition probability matrix $A = \{a_{ij}\}$
- Observation likelihoods
 - Output probability matrix $B = \{b_i(k)\}$
- Special initial probability vector π

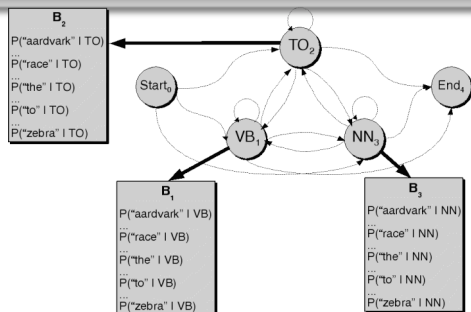
2/28/08

Transitions between the hidden states of HMM, showing A probs



2/28/08

B observation likelihoods for POS HMM



2/28/08

The A matrix for the POS HMM

	VB	TO	NN	PPSS
<s>	.019	.0043	.041	.067
VB	.0038	.035	.047	.0070
TO	.83	0	.00047	0
NN	.0040	.016	.087	.0045
PPSS	.23	.00079	.0012	.00014

Figure 4.15 Tag transition probabilities (the a array, $p(t_i|t_{i-1})$) computed from the 87-tag Brown corpus without smoothing. The rows are labeled with the conditioning event; thus $P(PPSS|VB)$ is .0070. The symbol <s> is the start-of-sentence symbol.

2/28/08

19

The B matrix for the POS HMM

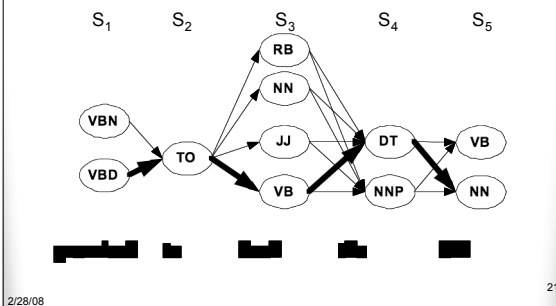
	I	want	to	race
VB	0	.0093	0	.00012
TO	0	0	.99	0
NN	0	.000054	0	.00057
PPSS	.37	0	0	0

Figure 4.16 Observation likelihoods (the b array) computed from the 87-tag Brown corpus without smoothing.

2/28/08

20

Viterbi intuition: we are looking for the best 'path'



2/28/08

21

The Viterbi Algorithm

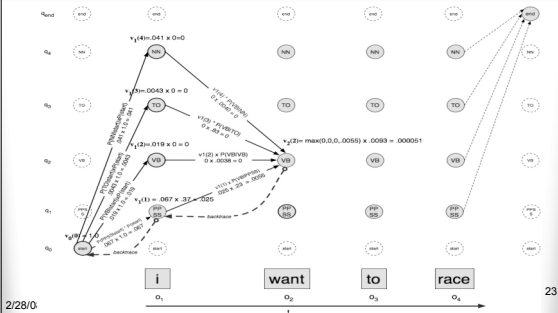
```

function VITERBI(observations of len  $T$ , state-graph of len  $N$ ) returns best-path
  create a path probability matrix  $viterbi[N+2, T]$ 
  for each state  $s$  from 1 to  $N$  do                : initialization step
     $viterbi[s, 1] \leftarrow a_{0,s} * b_s(o_1)$ 
     $backpointer[s, 1] \leftarrow 0$ 
  for each time step  $t$  from 2 to  $T$  do          : recursion step
    for each state  $s$  from 1 to  $N$  do
       $viterbi[s, t] \leftarrow \max_{s'} viterbi[s', t-1] * a_{s',s} * b_s(o_t)$ 
       $backpointer[s, t] \leftarrow \operatorname{argmax}_{s'} viterbi[s', t-1] * a_{s',s}$ 
     $viterbi[q_F, T] \leftarrow \max_{s=1}^N viterbi[s, T] * a_{s,q_F}$  ; termination step
     $backpointer[q_F, T] \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T] * a_{s,q_F}$  ; termination step
  return the backtrace path by following backpointers to states back in time from
   $backpointer[q_F, T]$ 
    
```

2/28/08

22

Viterbi example



2/28/0

23

Information Theory

- Who is going to win the World Series next year?
- Well there are 30 teams. Each has a chance, so there's a 1/30 chance for any team...? No.
 - ♦ Rockies? Big surprise, lots of information
 - ♦ Yankees? No surprise, not much information

2/28/08

24

Information Theory

- How much uncertainty is there when you don't know the outcome of some event (answer to some question)?
- How much information is to be gained by knowing the outcome of some event (answer to some question)?

2/28/08

25

Aside on logs

- Base doesn't matter. Unless I say otherwise, I mean base 2.
- Probabilities lie between 0 and 1. So log probabilities are negative and range from 0 ($\log 1$) to $-\infty$ ($\log 0$).
- The $-$ is a pain so at some point we'll make it go away by multiplying by -1 .

2/28/08

26

Entropy

- Let's start with a simple case, the probability of word sequences with a unigram model
- Example
 - ♦ $S = \text{"One fish two fish red fish blue fish"}$
 - ♦ $P(S) = P(\text{One})P(\text{fish})P(\text{two})P(\text{fish})P(\text{red})P(\text{fish})P(\text{blue})P(\text{fish})$
 - ♦ $\text{Log } P(S) = \text{Log } P(\text{One}) + \text{Log } P(\text{fish}) + \dots + \text{Log } P(\text{fish})$

2/28/08

27

Entropy cont.

- In general that's
- But note that
 - ♦ the order doesn't matter
 - ♦ that words can occur multiple times
 - ♦ and that they always contribute the same each time
 - ♦ so rearranging...



2/28/08

28

Entropy cont.

- One fish two fish red fish blue fish
- Fish fish fish fish one two red blue



2/28/08

29

Entropy cont.

- Now let's divide both sides by N, the length of the sequence:



- That's basically an average of the logprobs

2/28/08

30

Entropy

- Now assume the sequence is really really long.
- Moving the N into the summation you get

$$\frac{1}{N} \sum_{i=1}^N -\log_2 p_i$$

- Rewriting and getting rid of the minus sign

$$\sum_{i=1}^N \log_2 \frac{1}{p_i}$$

2/28/08

31

Entropy

- Think about this in terms of uncertainty or surprise.
 - ♦ The more likely a sequence is, the lower the entropy. Why?

$$\frac{1}{N} \sum_{i=1}^N -\log_2 p_i$$

2/28/08

32

Model Evaluation

- Remember the name of the game is to come up with statistical models that capture something useful in some body of text or speech.
- There are precisely a gazillion ways to do this
 - ♦ N-grams of various sizes
 - ♦ Smoothing
 - ♦ Backoff...

2/28/08

33

Model Evaluation

- Given a collection of text and a couple of models, how can we tell which model is best?
- Intuition... the model that assigns the highest probability to a set of withheld text
 - ♦ Withheld text? Text drawn from the same distribution (corpus), but not used in the creation of the model being evaluated.

2/28/08

34

Model Evaluation

- The more you're surprised at some event that actually happens, the worse your model was.
- We want models that minimize your surprise at observed outcomes.
- Given two models and some training data and some withheld test data... which is better?

2/28/08

35

Three HMM Problems

- Given a model and an observation sequence
 - ♦ Compute $\text{Argmax } P(\text{states} \mid \text{observation seq})$
 - Viterbi
 - ♦ Compute $P(\text{observation seq} \mid \text{model})$
 - Forward
 - ♦ Compute $P(\text{model} \mid \text{observation seq})$
 - EM (magic)

2/28/08

36

Viterbi

- Given a model and an observation sequence, what is the most likely state sequence
 - ♦ The state sequence is the set of labels assigned
 - ♦ So using Viterbi with an HMM solves the sequence classification task

2/28/08

37

Forward

- Given an HMM model and an observed sequence, what is the probability of that sequence?
 - $P(\text{sequence} \mid \text{Model})$
 - Sum of all the paths in the model that could have produced that sequence
 - So...
 - How do we change Viterbi to get Forward?

2/28/08

38

Who cares?

- Suppose I have two different HMM models extracted from some training data.
- And suppose I have a good-sized set of held-out data (not used to produce the above models).
- How can I tell which model is the better model?

2/28/08

39

Learning Models

- Now assume that you just have a single HMM model (π , A, and B tables)
- How can I produce a second model from that model?
 - ♦ Rejigger the numbers... (in such a way that the tables still function correctly)
 - ♦ Now how can I tell if I've made things better?

2/28/08

40

EM

- Given an HMM structure and a sequence, we can learn the best parameters for the model without explicit training data.
 - ♦ In the case of POS tagging all you need is unlabelled text.
 - ♦ Huh? Magic. We'll come back to this.

2/28/08

41

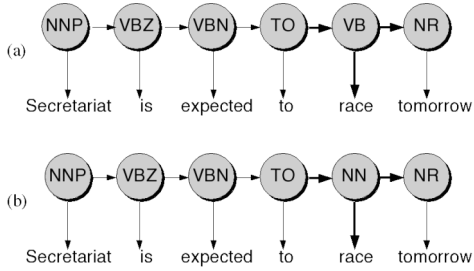
Generative vs. Discriminative Models

- For POS tagging we start with the question... $P(\text{tags} | \text{words})$ but we end up via Bayes at
 - ♦ $P(\text{words} | \text{tags})P(\text{tags})$
 - ♦ That's called a generative model
 - ♦ We're reasoning backwards from the models that could have produced such an output

2/28/08

42

Disambiguating "race"



2/28/08

43

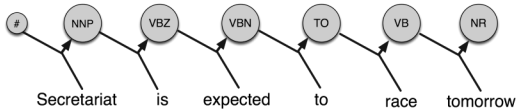
Discriminative Models

- What if we went back to the start to
 - ♦ $\text{Argmax } P(\text{tags}|\text{words})$ and didn't use Bayes?
 - ♦ Can we get a handle on this directly?
 - ♦ First let's generalize to $P(\text{tags}|\text{evidence})$
 - Let's make some independence assumptions and consider the previous state and the current word as the evidence. How does that look as a graphical model?

2/28/08

44

MaxEnt Tagging

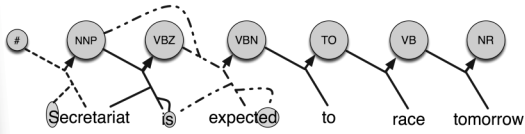


2/28/08

45

MaxEnt Tagging

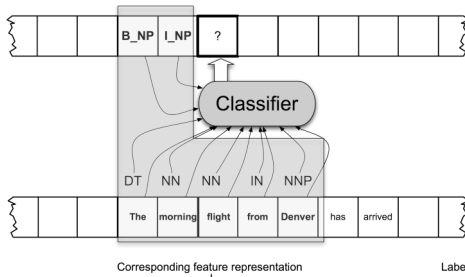
- This framework allows us to throw in a wide range of “features”. That is, evidence that can help with the tagging.



2/28/08

46

Statistical Sequence Classification



2/28/08

The, DT, B_NP, morning, NN, I_NP, flight, NN, from, IN, Denver, NNP, I_NP, has, arrived, I_NP

47
