

# **Part-of-Speech Tagging: Hidden Markov Models**

Katharina Kann — CSCI/LING5832

# Word Classes: Parts of Speech

- Traditional parts of speech:
  - Noun, verb, adjective, preposition, adverb, article, interjection, pronoun, conjunction, etc.
  - Lots of names for this notion: Part of speech, lexical category, word class, lexical tag...
  - Lots of debate within linguistics about the number, nature, and universality of these categories

# Word Classes: Parts of Speech

- Sources of evidence:
  - Morphological evidence
    - walk, walking, walked, walks
    - probably a verb

# Word Classes: Parts of Speech

- Sources of evidence:
  - Morphological evidence
    - walk, walking, walked, walks
    - probably a verb
  - Distributional evidence
    - The crash, A crash, Two crashes,  
The big crash...
    - probably a noun

# Penn Treebank POS Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb, gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, singular	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... --</i>
RP	particle	<i>up, off</i>			

# POS Tagging

- The process of assigning a part of speech or lexical class marker to each word in a text.
- Often a useful first step in an NLP pipeline.
- Fast and accurate taggers are widely available for many languages.
  - (Now we will learn how to make one!)

# POS Tagging

- The process of assigning a part of speech or lexical class marker to each word in a text.
- This is our first example of a **sequence labeling** task:
  - Assigning a category label to each element of a sequence.

# POS Tagging

WORD

tag

the

**DET**

koala

**N**

put

**V**

the

**DET**

keys

**N**

on

**P**

the

**DET**

table

**N**



# POS Tagging

- Many words have more than one part of speech:  
back
  - The **back** door = JJ
  - On my **back** = NN
  - Win the voters **back** = RB
  - Promised to **back** the bill = VB

# POS Tagging

- Many words have more than one part of speech:  
back
  - The **back** door = JJ
  - On my **back** = NN
  - Win the voters **back** = RB
  - Promised to **back** the bill = VB
- The POS tagging problem is to determine the tag for a particular instance of a word **in context**.
  - The context is usually a sentence.

# POS Tagging

- Note this is *distinct* from the task of word sense disambiguation.
  - “...**backed** the car into a pole”
  - “...**backed** the wrong candidate”

# Measuring Ambiguity

		87-tag Original Brown	45-tag Treebank Brown
<b>Unambiguous (1 tag)</b>		<b>44,019</b>	<b>38,857</b>
<b>Ambiguous (2–7 tags)</b>		<b>5,490</b>	<b>8844</b>
Details:	2 tags	4,967	6,731
	3 tags	411	1621
	4 tags	91	357
	5 tags	17	90
	6 tags	2 ( <i>well, beat</i> )	32
	7 tags	2 ( <i>still, down</i> )	6 ( <i>well, set, round, open, fit, down</i> )
	8 tags		4 ( <i>'s, half, back, a</i> )
	9 tags		3 ( <i>that, more, in</i> )

# Methods for POS Tagging

- Rule-based tagging
- Probabilistic sequence models
  - HMM (hidden Markov model) tagging
  - RNNs (recurrent neural networks)
  - [Transformers]

# Methods for POS Tagging

- Rule-based tagging
- Probabilistic sequence models
  - HMM (hidden Markov model) tagging
  - RNNs (recurrent neural networks)
  - [Transformers]
- Trivial baselines can often do well, too

# POS Tagging as Sequence Labeling

- We are given a sentence (an “observation” or “sequence of observations”).
  - I should book the flight
- What is the best sequence of tags that corresponds to this sequence of observations?
- Probabilistic view:
  - Consider all possible sequences of tags and assign a probability to each
  - Out of this universe of sequences, choose the tag sequence that is most probable given the observation sequence of our  $n$  words.

# Probabilistic Approach

- 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.

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



# Probabilistic Approach

- 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.

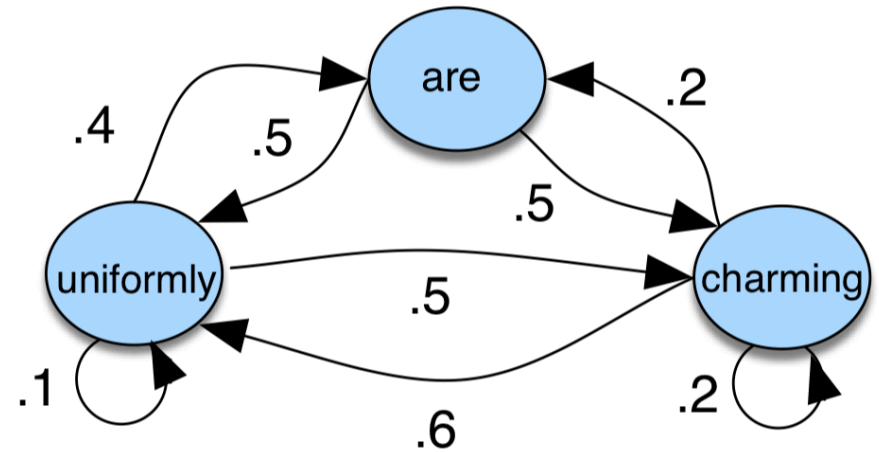
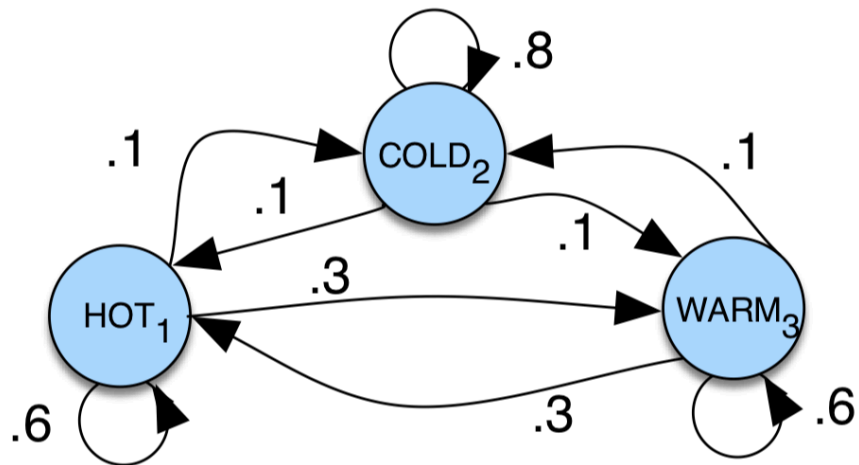
$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

# Markov Chains

# Markov Chain

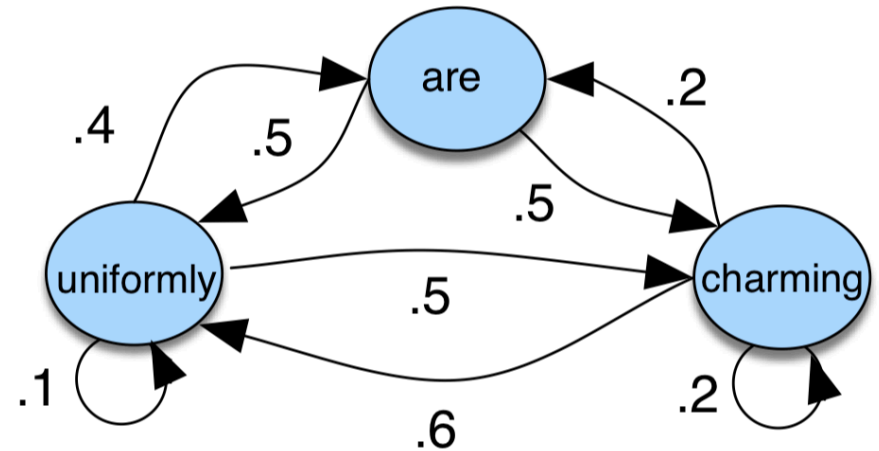
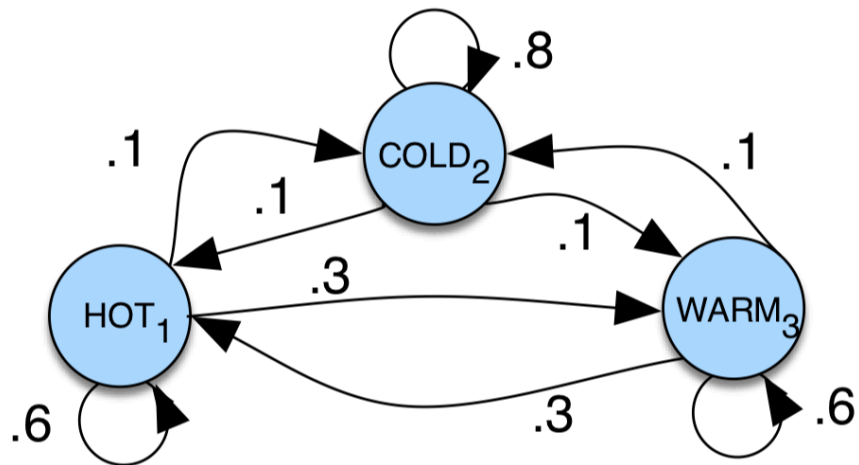
- A model that tells us something about the probabilities of sequences of random variables (*states*).
- Each of those variables can take on values from a predefined set.

# Markov Chain



J&M, Ch. 8

# Markov Chain



J&M, Ch. 8

- Markov assumption: only the last state matters!

# Markov Chain

$$\begin{aligned} P(q_1, \dots, q_T) = & P(q_T | q_1, \dots, q_{T-1}) \\ & * P(q_{T-1} | q_1, \dots, q_{T-2}) \\ & * \dots \\ & * P(q_1) \end{aligned}$$

**Markov assumption:**

$$P(q_T | q_1, \dots, q_{T-1}) = P(q_T | q_{T-1})$$

# Markov Chain

$$\begin{aligned} P(q_1, \dots, q_T) &= P(q_T | q_1, \dots, q_{T-1}) \\ &\quad * P(q_{T-1} | q_1, \dots, q_{T-2}) \\ &\quad * \dots \\ &\quad * \boxed{P(q_1)} \quad ? \end{aligned}$$

**Markov assumption:**

$$P(q_T | q_1, \dots, q_{T-1}) = P(q_T | q_{T-1})$$

# Markov Chain

- We further have an **initial probability distribution** over states.
  - The probability that each state is the first state!

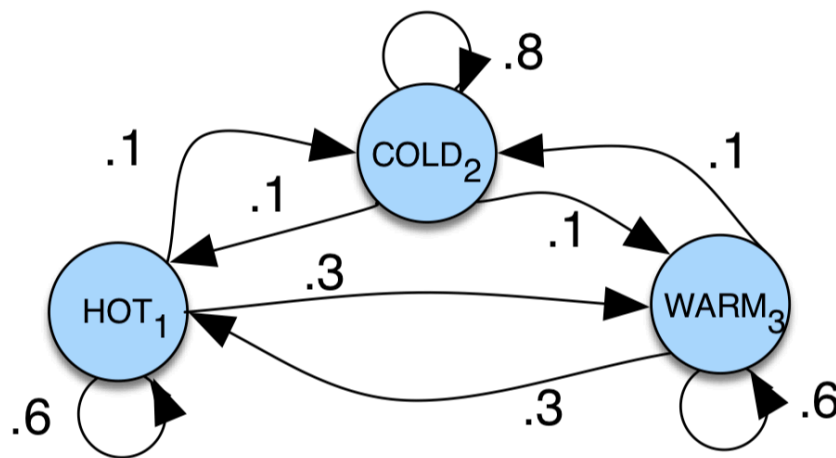


# Markov Chain

- We further have an **initial probability distribution** over states.
  - The probability that each state is the first state!
- For example:  $\pi = [0.1, 0.7, 0.2]$  for the states *hot*, *warm*, and *cold*

# In-Class Exercise

- $\pi = [0.1, 0.7, 0.2]$  for the states *hot*, *warm*, and *cold*



- **Compute the probability of:**
  1. Warm, hot, warm
  2. cold, hot, cold, hot

# Formal Definition: Markov Chain

$Q = q_1 q_2 \dots q_N$	a set of $N$ <b>states</b>
$A = a_{11} a_{12} \dots a_{n1} \dots a_{nn}$	a <b>transition probability matrix</b> $A$ , each $a_{ij}$ representing the probability of moving from state $i$ to state $j$ , s.t. $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$
$\pi = \pi_1, \pi_2, \dots, \pi_N$	an <b>initial probability distribution</b> over states. $\pi_i$ is the probability that the Markov chain will start in state $i$ . Some states $j$ may have $\pi_j = 0$ , meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$

J&M, Ch. 8

# Hidden Markov Models

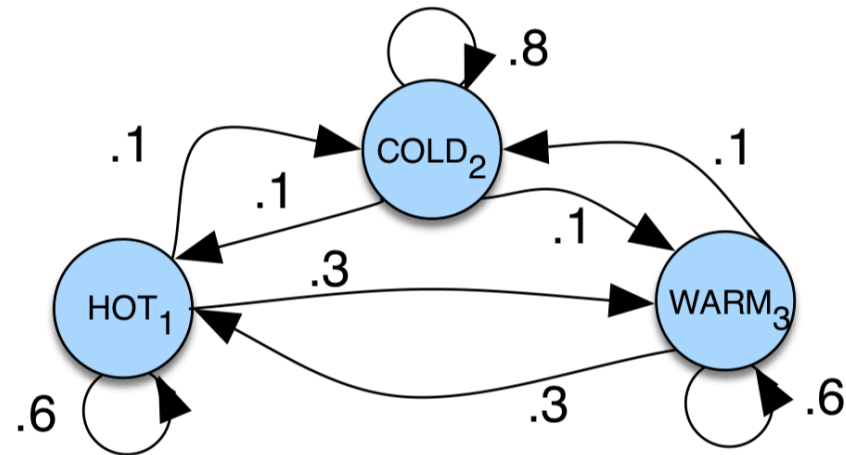
# Hidden Markov Model

- Markov chains can compute a probability for a sequence of *observable* events.
- However, we often don't observe the events we are interested in.
  - Hidden Markov models!

# Hidden Markov Model

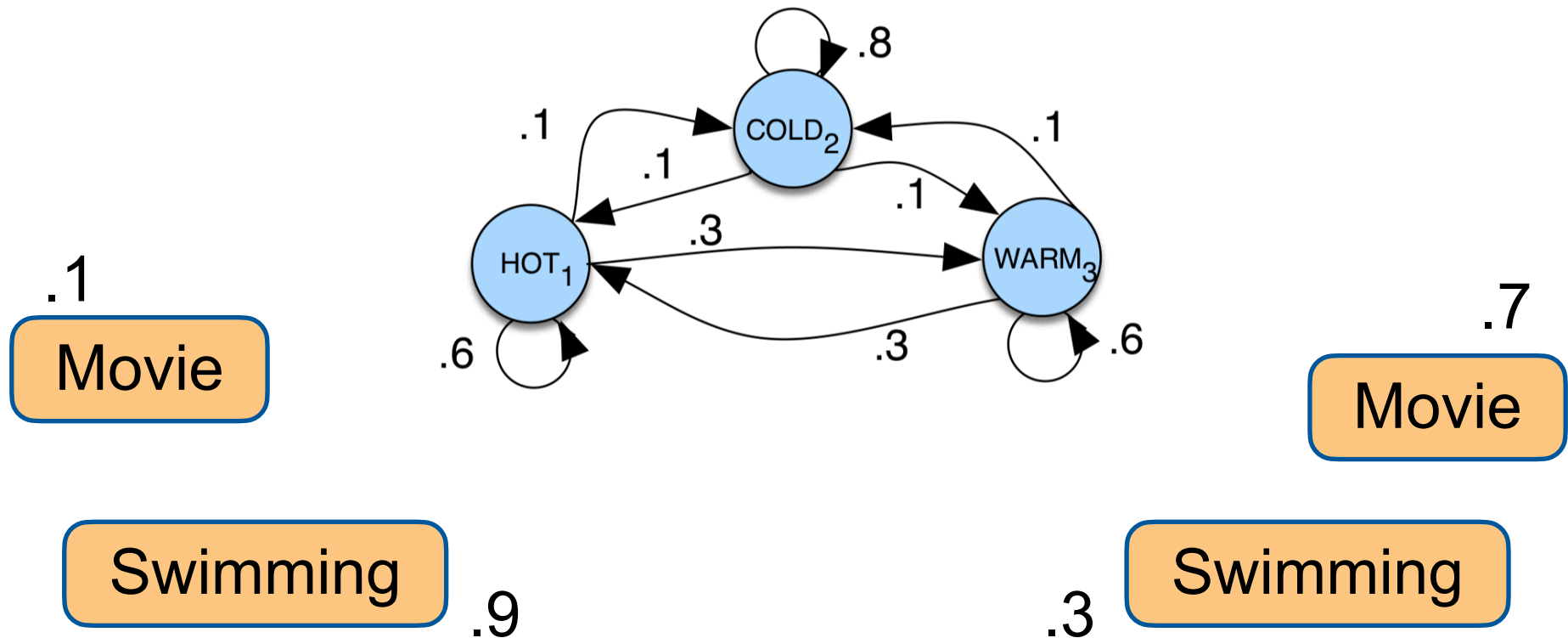
- Markov chains can compute a probability for a sequence of *observable* events.
- However, we often don't observe the events we are interested in.
  - Hidden Markov models!
- Example: part-of-speech tagging!

# Hidden Markov Model



# Hidden Markov Model

.2 Swimming Movie .8





# Hidden Markov Model

$Q = q_1 q_2 \dots q_N$	a set of $N$ <b>states</b>
$A = a_{11} \dots a_{ij} \dots a_{NN}$	a <b>transition probability matrix</b> $A$ , each $a_{ij}$ representing the probability of moving from state $i$ to state $j$ , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$
$O = o_1 o_2 \dots o_T$	a sequence of $T$ <b>observations</b> , each one drawn from a vocabulary $V = v_1, v_2, \dots, v_V$
$B = b_i(o_t)$	a sequence of <b>observation likelihoods</b> , also called <b>emission probabilities</b> , each expressing the probability of an observation $o_t$ being generated from a state $q_i$
$\pi = \pi_1, \pi_2, \dots, \pi_N$	an <b>initial probability distribution</b> over states. $\pi_i$ is the probability that the Markov chain will start in state $i$ . Some states $j$ may have $\pi_j = 0$ , meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$

# Hidden Markov Model

- A first-order Markov model makes 2 simplifying assumptions:

# Hidden Markov Model

- A first-order Markov model makes 2 simplifying assumptions:
  - Markov assumption:

$$P(q_i | q_1, \dots, q_{i-1}) = P(q_i | q_{i-1})$$

# Hidden Markov Model

- A first-order Markov model makes 2 simplifying assumptions:

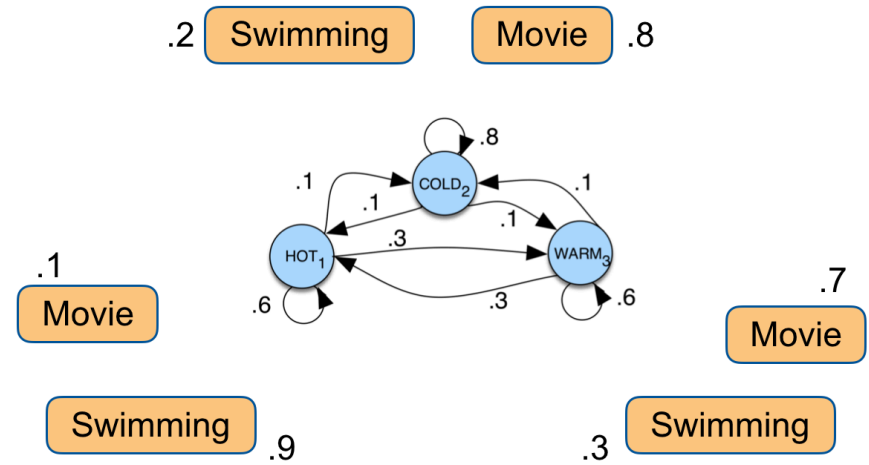
- Markov assumption:

$$P(q_i | q_1, \dots, q_{i-1}) = P(q_i | q_{i-1})$$

- Output independence:

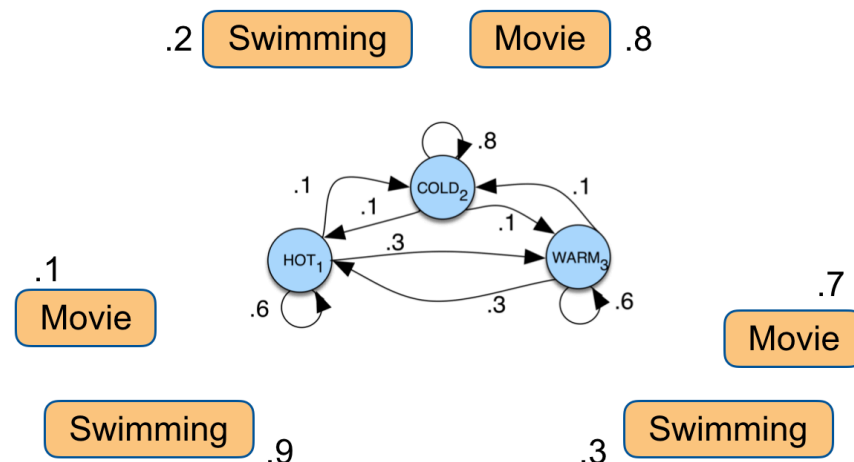
$$P(o_i | q_1, \dots, q_T, o_1, \dots, o_T) = P(o_i | q_i)$$

# Example



- Given:
  - $\pi = [0.1, 0.7, 0.2]$  for *hot*, *warm*, and *cold*
- Compute the probability of: [cold, warm, Movie, Movie]

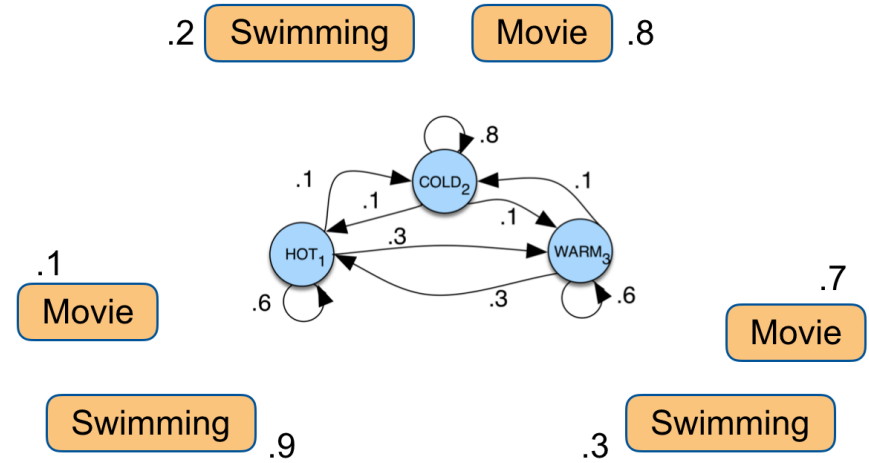
# Example



- Given:
  - $\pi = [0.1, 0.7, 0.2]$  for *hot*, *warm*, and *cold*
- Compute the probability of: [cold, warm, Movie, Movie]

$$P(c, w, M, M) = P(cold) * P(Movie | cold) * P(warm | cold) * P(Movie | warm)$$

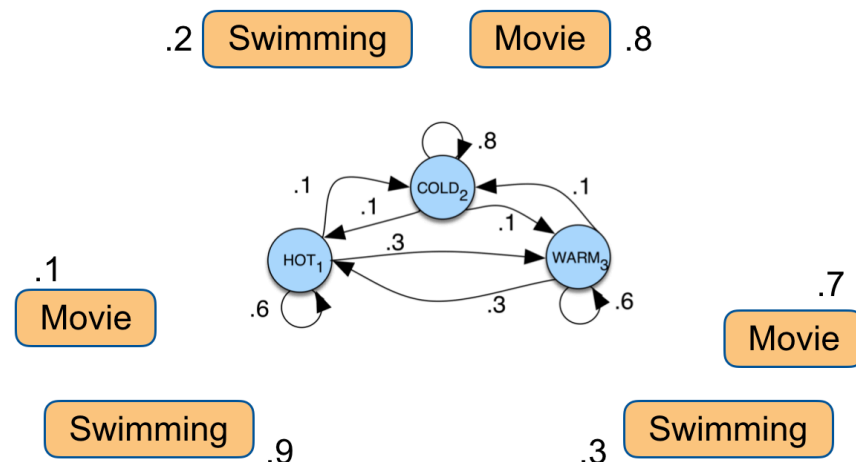
# Example



- Given:
  - $\pi = [0.1, 0.7, 0.2]$  for *hot*, *warm*, and *cold*
- Compute the probability of: [cold, warm, Movie, Movie]

$$\begin{aligned}
 P(c, w, M, M) &= P(cold) * P(Movie | cold) * \\
 &\quad P(warm | cold) * P(Movie | warm) \\
 &= 0.2 * 0.8 * \\
 &\quad 0.1 * 0.7
 \end{aligned}$$

# Example

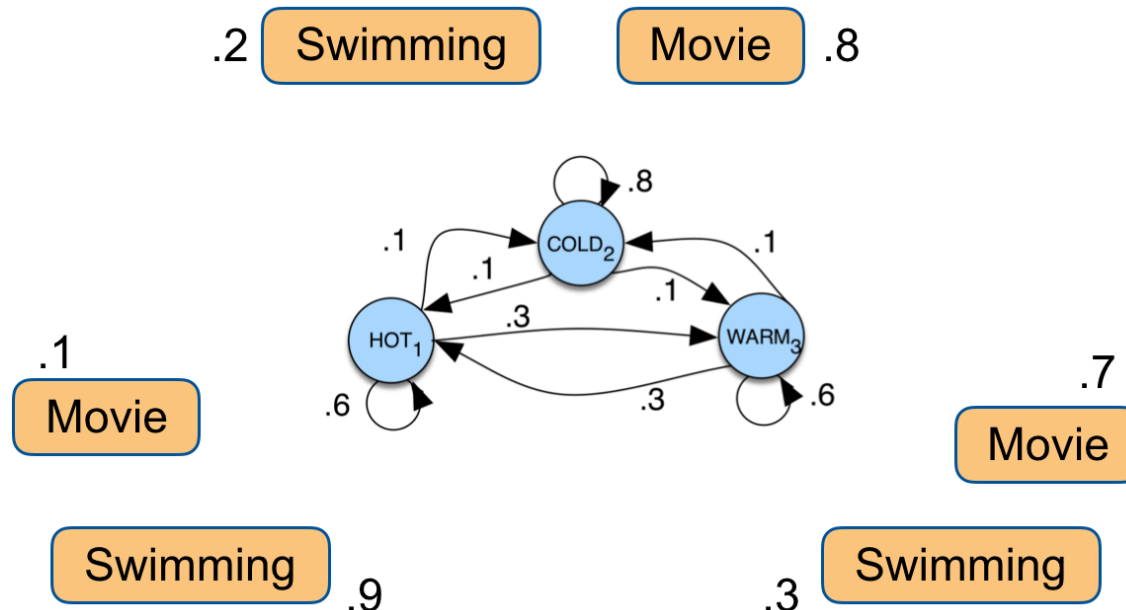


- Given:
  - $\pi = [0.1, 0.7, 0.2]$  for *hot*, *warm*, and *cold*
- Compute the probability of: [cold, warm, Movie, Movie]

$$\begin{aligned}
 P(c, w, M, M) &= P(cold) * P(Movie | cold) * \\
 &\quad P(warm | cold) * P(Movie | warm) \\
 &= 0.2 * 0.8 * \\
 &\quad 0.1 * 0.7 \\
 &= 0.0112
 \end{aligned}$$



# In-Class Exercise 2



- Given:  $\pi = [0.1, 0.7, 0.2]$  for the states *hot*, *warm*, and *cold*
- Compute the probability of:
  - [warm, cold, cold, Swimming, Movie, Movie]
  - [cold, hot, warm, Swimming, Movie, Movie]
  - [hot, cold, warm, Swimming, Movie, Movie]

# Back to Part-of-Speech Tagging

- 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  $P(\text{NN}|\text{DT})$  and  $P(\text{JJ}|\text{DT})$  to be high
- Compute  $P(\text{NN}|\text{DT})$  by counting in a labeled corpus:  
$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

# Back to Part-of-Speech Tagging

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

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

# Back to Part-of-Speech Tagging

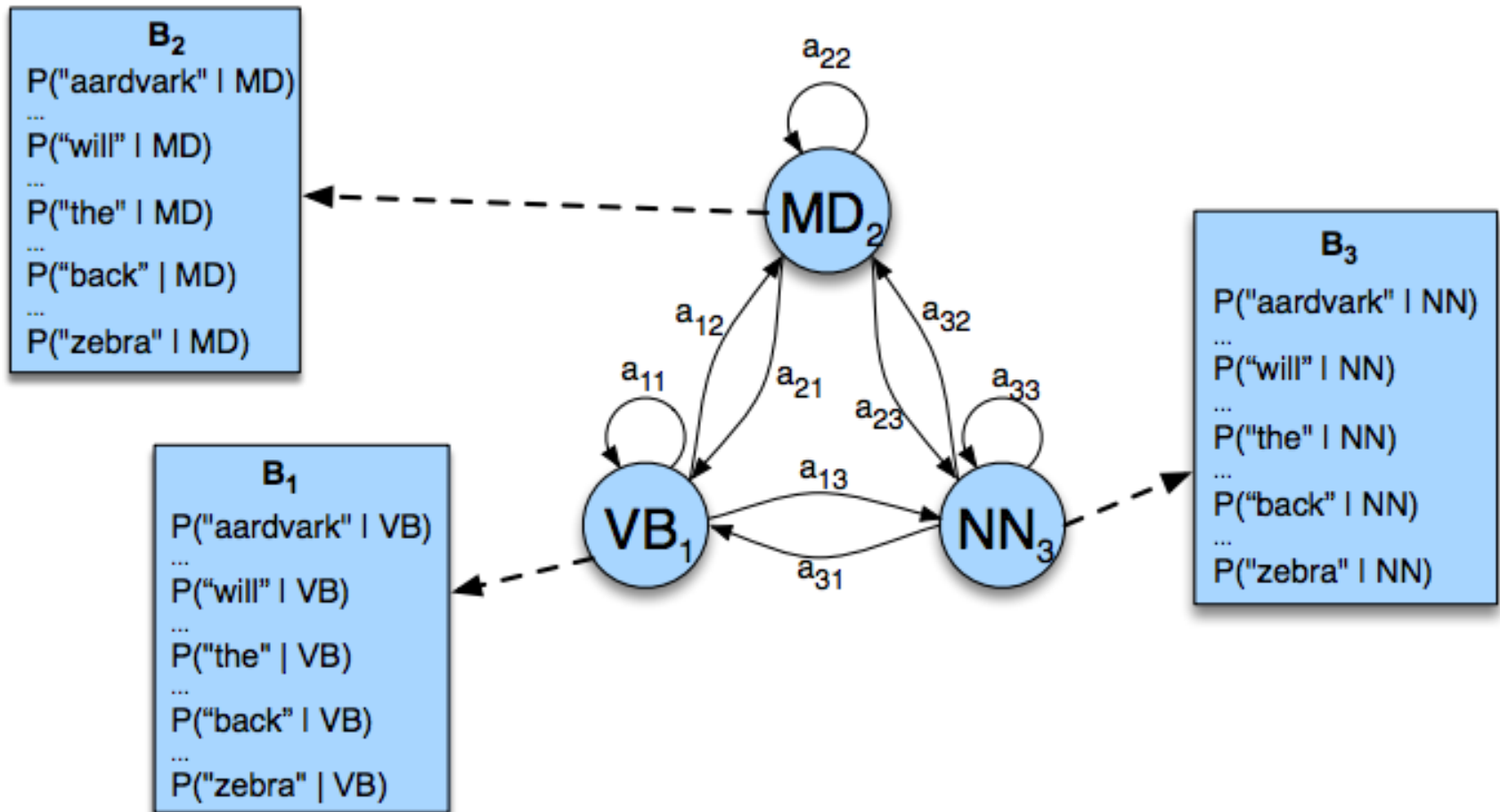
- Word likelihood probabilities  $p(w_i|t_i)$ 
  - VBZ (3sg Pres Verb) likely to be “is”
  - Compute  $P(\text{is}|\text{VBZ})$  by counting in a labeled corpus: 
$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

# Back to Part-of-Speech Tagging

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

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

# Back to Part-of-Speech Tagging



# 3 Problems

- Given this framework there are 3 problems that we can pose to an HMM:
  - Given an observation sequence and a model, what is the probability of that sequence?
  - Given an observation sequence and a model, what is the most likely state sequence?
  - Given an observation sequence, what are the best model parameters?

# 3 Problems

- Given this framework there are 3 problems that we can pose to an HMM:
  - Given an observation sequence and a model, what is the probability of that sequence?
  - Given an observation sequence and a model, what is the most likely state sequence?
  - Given an observation sequence, what are the best model parameters?

**We did a small part of this in the in-class exercise! (Why?)**



# Question

- If there are 30 or so tags in the Penn set...
- ...and the average sentence is around 20 words...
- ...how many tag sequences do we have to enumerate to argmax over?

# Question

- If there are 30 or so tags in the Penn set...
- ...and the average sentence is around 20 words...
- ...how many tag sequences do we have to enumerate to argmax over?

$$30^{20}$$

# Problem 1

- The probability of a sequence given a model...

**Computing Likelihood:** Given an HMM  $\lambda = (A, B)$  and an observation sequence  $O$ , determine the likelihood  $P(O|\lambda)$ .

- Used in sequence classification tasks
  - Word spotting in ASR, language identification, speaker identification, author identification, etc.

# Problem 1

- The probability of a sequence given a model...

**Computing Likelihood:** Given an HMM  $\lambda = (A, B)$  and an observation sequence  $O$ , determine the likelihood  $P(O|\lambda)$ .

- Used in sequence classification tasks
  - Word spotting in ASR, language identification, speaker identification, author identification, etc.
  - Train one model per class
  - Given an observation, pass it to each model and compute  $P(\text{seq}|\text{model})$
  - Argmax over models gives you the class

# Problem 2

- Most probable state sequence given a model and an observation sequence

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and a sequence of observations  $O = o_1, o_2, \dots, o_T$ , find the most probable sequence of states  $Q = q_1 q_2 q_3 \dots q_T$ .

# Problem 2

- Most probable state sequence given a model and an observation sequence

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and a sequence of observations  $O = o_1, o_2, \dots, o_T$ , find the most probable sequence of states  $Q = q_1 q_2 q_3 \dots q_T$ .

- Typically used in sequence labeling problems, where the labels correspond to hidden states
  - As we'll see almost any problem can be cast as a sequence labeling problem

# Problem 3

- Infer the best model parameters, given a model and an observation sequence...
  - That is, fill in the A and B tables with the right numbers...
    - The numbers that make the observation sequence most likely.

# Problem 3

- Infer the best model parameters, given a model and an observation sequence...
  - That is, fill in the A and B tables with the right numbers...
    - The numbers that make the observation sequence most likely.
  - Useful for getting an HMM without supervision/annotators!



# Solutions

- Problem 1: Forward
- Problem 2: Viterbi
- Problem 3: Forward-Backward
  - An instance of Expectation Maximization (EM)

# Wrapping up

- Discussed today:
  - Part-of-speech tagging
  - Markov chains
  - Hidden Markov models
- On Wednesday: The Viterbi algorithm  
(guest lecture by Abhidip)