

Faculty of Computer Science, Dalhousie University
CSCI 4152/6509 — Natural Language Processing

04-Nov-2025

Lecture 12: POS Tagging and Hidden Markov Model

Location: Studley LSC-Psychology P5260 Instructor: Vlado Keselj
 Time: 14:35 – 15:55

Previous Lecture

- N-gram model and Markov Chain model:
- Language modeling
- N-gram model assumption
- N-gram model graphical representation
- N-gram model as Markov chain
- Language model evaluation; Perplexity
- Text classification using language modeling
- N-gram Model Smoothing:
 - Add-one smoothing (Laplace smoothing)
 - *to continue*

13.1.2 Witten-Bell Discounting

In the context of data compression, Witten and Bell (1991) analyzed several smoothing methods, under the title “The Zero-Frequency Problem: Estimating the Probabilities of Novel Events in Adaptive Text Compression”. They considered three methods A, B, and C, and then, based on a Poisson process modelling, the methods P, X, and XC. It is interesting to note that the method X uses the same, or very similar, idea as in the Good-Turing smoothing.

The method C is what is usually referred to as the Witten-Bell smoothing. It uses an intuitive idea from data compression. Let us assume that we use a data compression method, which uses a dictionary of tokens w_1, w_2, \dots, w_r , so far. As long as the new tokens are from this set, we can encode them in some way, but whenever a new token appears, we need a special ‘escape’ code to introduce this token to the vocabulary. In this way, we can think of new tokens appearing as a new event, beside the events of seeing existing tokens. This is supported in practice by seeing approximately constant rate of new words appearing in a text after some initial reading. We can use the frequency of such ‘escape’ code as an estimate of the probability of seeing previously unseen events, and make sure that we allocate that much probability distribution mass for the smoothing purposes.

Example: Let us consider again using the example of training data ‘mississippi’ to train a unigram model, and then use it to estimate probability of the string ‘river’.

We will consider again estimating probability of 26 lower-case letters from the word ‘mississippi’. As in the case of unsmoothed n-grams, we will count letters: ‘m’ 1 time, ‘i’ 4 times, ‘s’ 4 times, and ‘p’ 2 times. However, we also note that we saw 4 different letters, which is equivalent to seeing ‘escape’ character 4 times, so we will reserve count 4 for unseen events in future as well. This is how we get the following probabilities using the Witten-Bell discounting:

Letter	Modified counts	Estimated frequency
i	4	$4/15 \approx 0.266666666666667$
m	1	$1/15 \approx 0.0666666666666667$
p	2	$2/15 \approx 0.133333333333333$
s	4	$4/15 \approx 0.266666666666667$
<i>new letters total</i>	4	$4/15 \approx 0.266666666666667$
Total:	15	

When we split the probability reserved for new letters equally among the remaining $26 - 4 = 22$ letters, we obtain the final estimated frequency:

Letter	Estimated frequency
i	$4/15 \approx 0.266666666666667$
m	$1/15 \approx 0.0666666666666667$
p	$2/15 \approx 0.133333333333333$
s	$4/15 \approx 0.266666666666667$
<i>other letters</i>	$\frac{4}{15 \cdot 22} = 2/165 \approx 0.0121212121212121$

The probability of the word 'river' in this model would be:

$$P(\text{'river'}) = P(r)P(i)P(v)P(e)P(r) = \frac{2}{165} \cdot \frac{4}{15} \cdot \frac{2}{165} \cdot \frac{2}{165} \cdot \frac{2}{165} \approx 5.75642615879697 \cdot 10^{-9}$$

Formulae for Witten-Bell Discounting If we want to express this in terms of formulae, we will denote that we saw r distinct tokens in a text of length n , or we can say that we saw n events and r 'escape' events, so the probability of seeing new tokens is $\frac{r}{n+r}$. Hence the unigram probability for seen tokens:

$$P(w) = \frac{\#(w)}{n+r}$$

and the total probability for unseen tokens is:

$$\frac{r}{n+r}$$

It is convenient that in the previous formulae we did not need to know the vocabulary size. If we do know the vocabulary size, we can now divide the probability for unseen tokens equally, and obtain:

$$P(w) = \frac{r}{(n+r)(|V| - r)}$$

for unseen tokens w .

Bigrams and Higher-order n-grams

The probabilities for bigrams and higher-order n-grams are smoothed in a similar way:

$$P(a|b) = \frac{\#(ba)}{\#(b) + r_b}$$

for seen bigrams ba , where r_b is the number of different tokens following b . The number $\#(b)$ does not represent necessarily the exact number of occurrences of b in this case. More precisely, it is the number of occurrences of b except at the end of text; i.e., the number of occurrences of b where b is followed by another token. The remaining probability mass for unseen events:

$$\frac{r_b}{\#(b) + r_b}$$

is used for unseen bigrams that start with b , and is usually divided according to lower-order n-grams; which would be unigrams in this case. If N_b is the set of all tokens that never follow b in the training data, then:

$$P(a|b) = \frac{r_b}{\#(b) + r_b} \cdot P(a) / \sum_{x \in N_b} P(x)$$

for unseen bigrams *ba*.

The next model: HMM. Our next probabilistic model is the Hidden Markov Model (HMM), and it is applicable to the task of labelling tokens of a sequence, such as the task of Part-of-Speech tagging (POS Tagging). Before that, we will make a review of the parts of speech in English, which are quite applicable with some changes to other natural languages as well.

Slide notes:

The Next Model: HMM

- HMM — Hidden Markov Model
- Typically used to annotate sequences of tokens
- Most common annotation: Part-of-Speech Tags (POS Tags)
- First, we will make a review of parts of speech in English

14 Part-of-Speech Tags (POS Tags)

Slide notes:

Part-of-Speech Tags (POS Tags)

- Reading: Sections 5.1–5.2 (Ch. 8 in new edition)
- Word classes called **Part-of-Speech (POS) classes**
 - also known as **syntactic categories, grammatical categories, or lexical categories**
- Ambiguous example: Time flies like an arrow.

	Time	flies	like	an	arrow.
1.	N	V	P	D	N
2.	N	N	V	D	N
	⋮				
- **POS tags:** labels to indicate POS class
- **POS tagging:** task of assigning POS tags

Note about reading material: Some reading material for the topics in this section can be found in the JM textbook in Sections 5.1–5.2 (“5.1 (Mostly) English Word Classes” and “5.2 Tagsets for English”), or in Chapter 8 of the upcoming edition 3 of the book.

The words in text are divided into classes according to their function, and these classes are called **Part-of-Speech classes** or **POS classes** for short. The POS classes are also sometimes called **syntactic categories, grammatical categories, or lexical categories**.

We can take a look at the ambiguous example of the sentence “Time flies like an arrow.” that we used before. The sentence can be interpreted in two ways, one is that the time goes very fast, and another one is that a species of flies, called “time flies”, like the arrow fruit. We can even think of a third meaning, which is a command to go and time the files immediately. If we label the words in this sentence according to their part-of-speech classes, we get three different sequences of labels for the the three different meanings as follows:

	Time	flies	like	an	arrow.
1.	N	V	P	D	N
2.	N	N	V	D	N
3.	V	N	P	D	N

The labels used below the words denote the following well-known POS classes: ‘N’ for nouns, ‘V’ for verbs, ‘P’ for prepositions, and ‘D’ for determiners.

The task of determining the part-of-speech label for each word in a sentence, or a text in general is called **POS**

tagging. From the above example, we can see that POS tagging is ambiguous, i.e., it may depend on the text interpretation by a reader.

POS Tag Sets

The concept of parts of speech as types of words used in language is known in linguistics for a long time. It was mentioned by several Antient Greek writers and it is well described in the work “The Art of Grammar”, which is believed to be written by Dionysius Thrax in the 2nd century BC. This work distinguishes the following eight parts of speech: nouns, verbs, pronouns, prepositions, adverbs, conjunctions, participles, and articles.

Slide notes:

POS Tag Sets

- Traditionally based on Ancient Greece source: eight parts of speech:
 - nouns, verbs, pronouns, prepositions, adverbs, conjunctions, participle, and articles
- Computer processing introduced a need for a large set of categories
- Useful in NLP, e.g.: named entity recognition, information extraction
- Various POS tag sets (in NLP):
 - Brown Corpus, Penn Treebank, CLAWS, C5, C7, ...
- We will use the Penn Treebank system of tags

WSJ Dataset

- WSJ — Wall Street Journal data set
- Most commonly used to train and test POS taggers
- Consists of 25 sections, about 1.2 million words
- Example:

Pierre NNP Vinken NNP , , 61 CD years NNS old JJ , ,
 will MD join VB the DT board NN as IN a DT
 nonexecutive JJ director NN Nov. NNP 29 CD . .
 Mr. NNP Vinken NNP is VBZ chairman NN of IN
 Elsevier NNP N.V. NNP , , the DT Dutch NNP
 publishing VBG group NN . .
 Rudolph NNP Agnew NNP , , 55 CD years NNS old JJ
 and CC former JJ chairman NN of IN Consolidated NNP
 Gold NNP Fields NNP PLC NNP , , was VBD named VBN

Open and Closed Categories

- Word POS categories are divided into two sets: *open* and *closed* categories:
- **open categories**
 - dynamic set
 - content words
 - larger set
 - e.g.: nouns, verbs, adjectives
- **closed categories** or **functional categories**:
 - fixed set

- small set
- frequent words
- e.g.: articles, auxiliaries, prepositions

The words of a language, and generally POS categories, can be divided into two sets: *open and closed categories*.

Open POS categories are lexical categories that are dynamic in sense that their content changes over time, or depending on dialect or domain of usage. These sets are larger and the words bear most of the information content in a text. Some examples of open word categories are nouns, verbs, and adjectives.

Closed or functional POS categories are lexical categories consisting of fixed sets of words, which are used frequently in text and they are typically used in a functional way, i.e., as syntactic fillers and with less information content. Examples of such categories are articles, auxiliaries, and prepositions.

14.1 Open Word Categories

The open word categories are: *nouns*, *adjectives*, *verbs*, and *adverbs*. There are also groups of adverbs that belong to the closed word categories. Generally, like many other rules in natural language, this division is not strict and there are many exceptions.

Open Word Categories

- nouns (NN, NNS, NNP, NNPS)
 - concepts, objects, people, and similar
- adjectives (JJ, JJR, JJS)
 - modify (describe) nouns
- verbs (VB, VBP, VBZ, VBG, VBD, VBN)
 - actions
- adverbs (RB, RBR, RBS)
 - modify verbs, but other words too

Nouns (NN, NNS, NNP, NNPS)

Nouns refer to people, animals, objects, concepts, and similar.

Features:

- number: singular, plural
- case: subject (nominative), object (accusative)
- Some languages have more cases, and more number values
- Some languages have grammatical gender

Nouns have a number of linguistic properties called *features*, which vary across languages. The features are important in creating larger phrases and sentences in a linguistically correct way. The most common feature is *number* and the main values for it are *singular* and *plural*, expressing whether we are talking about one instance of an object or multiple instances. Some languages distinguish a more finer-grained set of values for number. The number, as other features, is typically expressed by modifying a suffix of a word; for example, the suffix *-s* is added in English for plural nouns.

Case is another common feature for nouns, which is not used much in English. The case indicates the syntactic and semantic role a noun plays in a phrase or a sentence. For example, the *nominative* case is used for nouns in the subject position in a sentence, while *accusative* case is used in the direct object position.

Noun Tags and Examples

Slide notes:

Noun Tags and Examples	
NN	for common singular nouns; e.g., company, year, market
NNS	for common plural nouns; e.g., shares, years, sales, prices, companies
NNP	for proper nouns (names); e.g., Bush, Japan, Federal, New York, Corp, Mr., Friday, James A. Talcott (“James NNP A. NNP Talcott NNP”)
NNPS	for proper plural nouns; e.g., Canadians, Americans, Securities, Systems, Soviets, Democrats

The noun tags in the Penn tag set are: NN, NNS, NNP, and NNPS.

NN is used for common singular nouns, such as *company*, *year*, and *market*.

NNS is used for common plural nouns, such as *shares*, *years*, *sales*, *prices*, and *companies*.

NNP is used for proper singular nouns (names), which are the names of people, geographical entities, countries, institutions, and similar, such as *Bush*, *Japan*, *Federal*, *New York*, *Corp*, *Mr.*, *Friday*, *James A. Talcott*. The proper nouns consisting of several words are all tagged with the NNP tag; for example: “James NNP A. NNP Talcott NNP” The token “Mr.” comes with a person’s name and it would be tagged as NNP as well. The words “Federal” and “Corp.” are proper nouns as parts of institution or organization names. Somewhat specific to English, the names of days in a week and months, such as “Friday” and “January” are also proper nouns, so tagged as NNP.

NNPS is used for proper plural nouns, such as *Canadians*, *Americans*, *Securities*, *Systems*, *Soviets*, and *Democrats*.

Adjectives (JJ, JJR, JJS)

Adjectives describe properties of nouns; for example: red rose, long journey, etc.

Adjectives have three forms and each of them is tagged separately:

Form	Example	Tag
positive	rich	JJ
comparative	richer	JJR
superlative	richest	JJS

In the Brown corpus, the corresponding tags were JJ, JJR, and JJT; while the JJS tag was reserved for the **semantic superlative forms**, such as: chief, main, top, etc. These forms are tagged as JJ in the Penn Treebank corpus.

Comparatives and superlatives of longer adjectives in English are formed as multi-word sequences, such as “more intelligent” and “the most intelligent” for the adjective “intelligent”. These sequences are called the **periphrastic adjective forms**. They are tagged as follows:

```
more JJR intelligent JJ
and
the DT most JJS intelligent JJ
```

Verbs (VB, VBP, VBZ, VBG, VBD, VBN)

Verbs are used to describe:

- actions; e.g., throw the stone
- activities; e.g., walked along the river
- or states; e.g., have \$50

Verbs can have different forms and they are tagged accordingly:

Tag	Form name	Example
VB	base	eat, be, have, walk, do
VBD	past	ate, said, was, were, had
VBG	present participle	eating, including, according, being
VBN	past participle	eaten, been, expected
VBP	present non-3sg	eat, are, have, do, say, 're, 'm
VBZ	present 3sg	eats, is, has, 's, says

Gerund is a noun which has the same form as the present participle; e.g., 'Walking is fun.'

Verb Features:

- number: singular, plural
- person: 1st, 2nd, 3rd
- tense: present, past, future
- aspect: progressive, perfect
- mood: possibility, subjunctive (e.g. 'They requested that he be banned from driving.')
- participles: present participle, past participle
- voice: passive, active: "He wrote a book." vs. "A book was written by him."

Verb Tenses:

- present: I walk
- infinitive: to walk
- progressive: I am walking
- present perfect I have walked
- past perfect: I have walked

Adverbs (RB, RBR, RBS)

Adverbs modify verbs, as their name suggests, but also other lexical classes, such as adjectives and adverbs. In this way their function is quite heterogeneous. For example, some typical adverbs that can be used to modify verbs are *allegedly* and *quickly*, because they obviously describe whether something happened or how something happened.

Not all adverbs belong to the open group of categories: a group of adverbs called **qualifiers** or **degree adverbs** belong to the closed group. Example of such adverbs are *very* and *not*.

Here are a few examples of adverb usage. An example of an adverbs modifying a verb:

She *often* travels to Las Vegas.

an example of adverbs modifying verbs and adverbs:

Unfortunately, John walked *home extremely slowly yesterday*.

and two examples of adverbs modifying adjectives:

a *very* unlikely event
a *shockingly* frank exchange

Adverb Inflections

Adverbs can have three forms, similarly to adjectives;

Tag	Form	Examples
RB	positive	late, often, quickly
RBR	comparative	later, better, less
RBS	superlative	most, best

The superlative adverbs are tagged as RBT in the Brown corpus.

Adverbial Nouns are nouns that also behave as adverbs. Such nouns are ‘home’ and ‘tomorrow.’ For example, we can say

I am going home.

but not

* I am going room.

In the Brown corpus these nouns were tagged as NNR, but in the Penn Treebank corpus they are tagged as NN.

Another noun tag in the Brown corpus that cannot be found in the Penn Treebank corpus is NN\$, which was used to denote possessives, like ‘people’s’; in the Penn Treebank this would be tagged as ‘people NNP ’s POS’.

14.2 Closed Word Categories

- small, fixed, frequent, functional group
- typically no morphological transformations
- include:
 - determiners, pronouns, prepositions, particles, auxiliaries and modal verbs, qualifiers, conjunctions, numbers, interjections

Determiners (DT)

- articles: the, a, an
- demonstratives:
 - this, that, those
 - some, any
 - either, neither
- quantifiers: all, some

Interrogative Determiners (WDT)

Interrogative determiners are tagged as a separate class. Some examples are: ‘what’, ‘which’, ‘whatever’, and ‘whichever’.

Predeterminers (PDT)

- Examples: both, quite, many, all such, half
- Examples in context:
“half the debt”, “all the negative campaign”
- Interesting classifications of determiners (Bond 2001)
 - by linear order: pre-determiners, central determiners, post-determiners
 - by meaning: quantifiers, possessives, determinatives

A Side Note:

Two interesting classifications of determiners were given by Francis Bond in his dissertation “Determiners and Numbers in English contrasted with Japanese, as exemplified in Machine Translation.” These classifications are a classification of determiners by linear order, and a classification by meaning. This classification is not in accordance with the Penn tag set; e.g., the numerals are also included in the set of determiners.

Determiners grouped by linear order:

- pre-determiners
 - quantifiers: *all, both*
 - fractions: *half, one-third, ...*
 - multiples: *double, twice, three times, ...*
 - *what* (exclamative: *What a great party!*)
- central determiners
 - articles: *the, a(an), some, any*
 - demonstratives: *this/these, that/those*
 - possessive pronouns: *my, your, his, her, its, their, our*
 - possessive phrases: *the king's, his friend's, ...*
 - quantifiers *no, some, any, either, neither, another, each, enough, much, more, most, less, a few, a little, many a, several*
 - *which, what* (interrogative: *What sound is that?*)
 - pronouns: *we, us, you*
- post-determiners
 - cardinal numerals: *one, two, three, ...*
 - fixed-numbers: *dozen, score, ...*
 - quantifiers: *every, many, few, little*
 - emphatic possessive: *own*
 - *such*

Determiners grouped by meaning:

- quantifiers
 - cardinal numerals *one, two, three, ...*
 - other quantifiers *all, both, no, some, any, much, many, few, a few, little, a little, either, neither, another, enough, more, most, less, many a, several*
 - fractions: *half, one-third, ...*
 - multiples: *double, twice, three times, ...*
- possessives
 - possessive pronouns: *my, your, his, her, its, their, our*
 - possessive phrases: *the king's, his friend's, ...*
 - emphatic possessive: *own*
- determinatives
 - articles: *the, a/an, some/any*
 - demonstratives: *this/these, that/those*

- *which, what* (interrogative)
- *what* (exclamative)
- *such*
- pronouns: *we, us, you*

Pronouns (PRP, PRP\$)

- PRP for personal pronouns
 - examples: I, you, he, she, it, we, you, they
- PRP tag for accusative case (diff. tag in Brown):
 - examples: me, him, her, us, them
- PRP tag for reflexive pronouns (diff. in Brown):
 - examples: myself, ourselves, ...
- PRP\$ tag for possessive pronouns:
 - examples: your, my, her, his, our, their, its
- PRP for second possessives (diff. in Brown):
 - examples: ours, mine, yours, ...

The personal pronouns are tagged with PRP in the Penn tagset. The following are some of the features of pronouns:

- number: singular (sg), plural (pl)
- person: first (1st), second (2nd), third (3rd)
- case: nominative (subject), accusative (object)
- gender: masculine (he), feminine (she), neuter (it)

The singular personal pronouns used to be tagged with PP in the Brown corpus, and the plural personal pronouns were tagged with PPS (we, you, they).

The personal pronouns in accusative case (me, you, him, her, it, us, you, them) have the same PRP tag, while in the Brown corpus they had tag PPO. The reflexive pronouns (myself, ourselves, ...) have the same tag PRP, while they used to be tagged PPL and PPLS in Brown.

The tag for **possessive pronouns** is PRP\$; e.g., for your, my, her, our, his, their, its.

The **second possessives** (ours, mine, yours, ...) are tagged PRP (they used to be tagged PP\$\$ in Brown).

Wh-pronouns (WP) and Wh-possessive (WP\$)

- wh-pronouns (WP): who, what, whom, whoever, ...
- wh-possessive pronoun (WP\$): whose

Prepositions (IN)

Prepositions reflect spatial or time relationships.

Examples: of, in, for, on, at, by, concerning, ...

Particles (RP)

- frequently ambiguous and confused with prepositions
- used to create compound verbs
- examples: put off, take off, give in, take on, “went on for days,” “put it off”

Possessive ending (POS)

- possessive clitic: 's
- Example: John's book
- tagged as: John NNP 's POS book NN

Modal Verbs (MD)

- the examples of modal verbs: can, may, could, might, should, will
- and their abbreviations: 'd, 'll
- tag for modal verbs: MD
- negative forms are separated into a modal verb and an adverb 'not' (will be covered); e.g.: 'couldn't' is tagged as "could MD n't RB"
- *Auxiliary verbs* are: be, have, and do; and their different forms
- in Brown: each auxiliary verb has a separate tag
- in Penn Treebank: they are tagged in the same way as common verbs (we will see that later)

Infinitive word 'to' (TO)

- used to denote an infinitive: e.g., to call
- 'na' is marked as TO in 'gonna', 'wanna' and similar; e.g.: "gonna call" is tagged "gon VB na TO call VB"

Qualifiers (RB)

- qualifiers are closed adverbs, and they are tagged as adverbs (RB) (covered later)
- example: not, n't, very
- postqualifiers: enough, indeed

Wh-adverbs (WRB)

Examples: how, when, where, whichever, whenever,...

Conjunctions (CC)

- words that connect phrases
- coordinate conjunctions (tag: CC) connect coordinate phrases:
- examples; and, or, but, yet, plus, versus, ...
- subordinate conjunctions connect phrases where one is subordinate to another
- examples: if, although, that, because, ...
- subordinate conjunctions are tagged as prepositions (IN) in Penn Treebank
- in Brown corpus, they used to be tagged CS

Numbers (CD)

Numbers behave in a similar way to adjectives: they also modify nouns. Here, we distinguish two kinds of numbers: **cardinal numbers** or **cardinals**, and **ordinal numbers** or **ordinals**.

Examples:

- cardinals: 1, 0, 100.34, hundred, etc.
- ordinals: first, second, 3rd, 4th, etc.

Cardinal numbers are tagged as CD.

Ordinal numbers have a separate tag in the Brown corpus—OD. In the Penn Treebank corpus, they are tagged as adjectives—JJ.

Interjections (UH)

Examples: yes, no, well, oh, quack, OK, please, indeed, hello, Congratulations, . . .

14.3 Remaining POS Classes

Existential ‘there’ (EX)

Belongs to closed word category; example: “There/EX are/VBP three/CD classes/NNS per/IN week/NN”

Foreign Words (FW)

Examples: de (tour de France), perestroika, pro, des

List Items (LS)

Examples: 1, 2, 3, 4, a., b., c., first, second, etc.

Punctuation

Examples	Tag	Description
,	,	comma
; : . . . - --	:	mid-sentence separator
. ! ?	.	sentence end
({ [<	(open parenthesis
) }] >)	closed parenthesis
` `` non-``	``	open quote
' ''	''	closed quote
\$ c HK\$ CAN\$	\$	dollar sign
#	#	pound sign
- + & @ * ** ffr	SYM	everything else

14.4 Overview of POS Tags

Penn Treebank Part-of-Speech Tags (adapted from [JM])

Tag	Description	example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two, three</i>	TO	infinitive 'to'	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjections	<i>uh, oops</i>
EX	existential 'there'	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>it mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition or subordinating conjunction	<i>of, in, by</i>	VBG	verb, present participle	<i>eating</i>
JJ	adjective	<i>rich</i>	VBN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>richer</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>richest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, a</i>	WDT	wh-determiner	<i>which</i>
MD	modal verb	<i>can should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, singular or mass	<i>llama, snow</i>	WP\$	wh-possessive	<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>Canadians</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	"	left quote	<i>' , "</i>
POS	possessive ending	<i>'s</i>	"	right quote	<i>' , "</i>
PRP	personal pronoun	<i>I, you, we</i>	(left parenthesis	<i>(, [</i>
PRP\$	possessive pronoun	<i>your, one's</i>)	right parenthesis	<i>),]</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-end punc.	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc.	<i>: ;</i>
RP	particle	<i>up, off</i>			

14.5 Some Tagged Examples

The/DT grand/JJ jury/NN commented/VBD on/IN
a/DT number/NN of/IN other/JJ topics/NNS ./.

Book/VB that/DT flight/NN ./.

Does/VBZ that/DT flight/NN serve/VB dinner/NN ?/.

It/PRP does/VBZ a/DT first-rate/JJ job/NN ./.

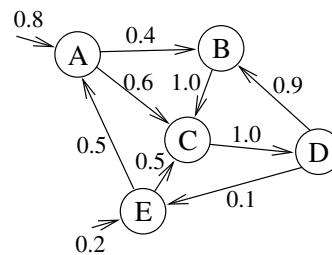
``/`` When/WRB the/DT sell/NN programs/NNS hit/VBP
,/, you/PRP can/MD hear/VB the/DT order/NN
printers/NNS start/VB to/TO go/VB ''/'' on/IN the/DT
Big/NNP Board/NNP trading/NN floor/NN ,/, says/VBZ
one/CD specialist/NN there/RB ./.

``/`` Do/VBP you/PRP make/VB sweatshirts/NNS or/CC
sparkplugs/NNS ?/.

15 Hidden Markov Model (HMM)

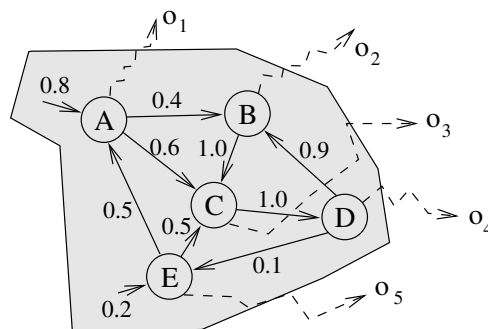
Let us take consider the problem of part-of-speech tagging; i.e., POS tagging. A POS tagger would be an algorithm that takes a tokenized sentence as an input and produces a tagged sentence as the output; i.e., the same sentence in which each token is associated with one of the POS tags. If we want to solve this problem using probabilistic modeling, then it is natural to associate all tokens to probabilistic variables, and their tags as well. There are dependencies between words and their associated tags, but it also seems that tags form some typical sequences, so there are dependencies from each tag to the following tag. This is a motivation for introducing our next model, the Hidden Markov Model.

We will look again at the example of Markov Chain, shown in a section before.



Markov Chain Example

If we assume that the states in such model are not observable, i.e., that they are “hidden,” and we can actually observe only an “emitted” symbol, based on a probabilistic distribution for producing observable symbols given a hidden state, we obtain the *Hidden Markov Model* (HMM). An example of such model is represented in the following figure:



Hidden Markov Model Example

15.1 HMM Formal Definition

Slide notes:

HMM Formal Definition

- Five-tuple: (Q, π, a, V, b) (there are other variations)
- 1. set of states $Q = \{q_1, q_2, \dots, q_N\}$
- 2. initial distribution $\pi: \pi(q)$ for each state q
- 3. transition probabilities $a: a(q, s)$ for any two states q and s
- 4. output vocabulary $V = \{o_1, o_2, \dots, o_m\}$
- 5. output probability $b: b(q, o)$ for each state q and observable o

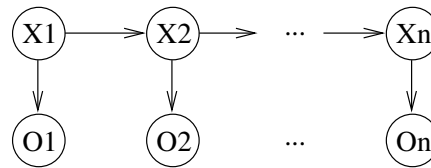
In more precise terms, an HMM (Hidden Markov Model) includes a finite set of hidden states $Q = \{q_1, q_2, \dots, q_N\}$. A probability distribution $\pi(q)$ ($0 \leq \pi(q) \leq 1$) specifies for each state q the probability that it will be the initial state, where the following constraint $\sum_q \pi(q) = 1$ holds. Instead of making transitions from state to state in a deterministic way, for each pair of states q_i and q_j there is a probability of the transition from state q_i to q_j $a_{ij} = a(q_i, q_j)$, so that $0 \leq a_{ij} \leq 1$ for each parameter a_{ij} and $\sum_s a(q, s) = 1$ for each state q . From each visited state an observable o is generated, where $o \in V$, and V is a finite vocabulary. For an arbitrary state $q \in Q$ and any observable symbol $o \in V$, the output probability $b(q, o)$ that the observable symbol o will be generated is defined, so that for all q $\sum_o b(q, o) = 1$.

We can summarize this into the following definition of a Hidden Markov Model:

Definition 15.1 (Hidden Markov Model) A Hidden Markov Model is a five-tuple (Q, π, a, V, b) , where: Q is a finite set of states $Q = \{q_1, q_2, \dots, q_N\}$, and $N \geq 1$; $\pi : Q \rightarrow [0, 1]$ is the initial probability distribution for the first state, with the constraint $\sum_{q \in Q} \pi(q) = 1$; $a : Q \times Q \rightarrow [0, 1]$ is the transition probability, $a(q_i, q_j)$ is also denoted as a_{ij} and it denotes probability of the next state q_j given the current state q_i , so $\sum_{q \in Q} a(s, q) = 1$; V is a finite output vocabulary $V = \{o_1, o_2, \dots, o_m\}$; and $b : Q \times V \rightarrow [0, 1]$ is the output probability, so that $b(q, o)$ is probability of generating the symbol o in the state q , $\sum_{o \in V} b(q, o) = 1$, and we also denote $b(q_i, o_j)$ as b_{ij} .

HMM Assumption

Given an HMM, we can generate samples by generating an initial state, producing an observable corresponding to that state, and then creating the next state, another observable produced by this state, and so on. For a particular length n , the following graph can be used to illustrate operation of an HMM:



This representation is sometimes called *unrolled* HMM graphical representation, in particular when compared with the DFA-style representation that we saw before. This unrolled representation is similar to the previous graphical representation of the Naïve Bayes model, and it is called the Belief Network, or Bayesian Network representation. We will later introduce a more general concept of the Bayesian network.

The value of X_1 is the initial state of the HMM, and the value of each consecutive variable X_i is the consecutive state of HMM. The values of variables O_1, O_2, \dots , are produced observables.

The HMM assumption formula is:

$$P(X_1, O_1, \dots, X_n, O_n) = P(X_1) \cdot P(O_1|X_1) \cdot P(X_2|X_1) \cdot P(O_2|X_2) \cdot \dots \cdot P(X_n|X_{n-1}) \cdot P(O_n|X_n)$$

HMM Application Areas

- Language Modelling
- Acoustic Modelling
- Part-of-Speech tagging (POS tagging)
- Many kinds of sequence tagging (e.g., extracting bio-medical terms)

HMMs are successfully used for acoustic modeling in speech recognition. They are also successfully applied to language modeling and POS tagging, among many applications in NLP.

15.2 POS Tagging using HMM

We will examine now Hidden Markov Model (HMM) in more details, including computational tasks in this model on the example of POS tagging application.

HMM use in POS Tagging

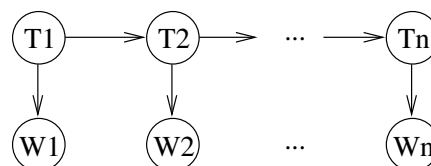
- Hidden states = POS Tags
- Observable variables = words
- In practice: higher-order HMM taggers are used, where the nodes keep a bit longer history (e.g., two previous tags)
- Described in [JM] Sec 5.5 (HMM POS Tagging)

Computational Tasks for HMM

- Evaluation: use HMM assumption formula
- Generation: generate in the order dictated by the “unrolled” graphical representation
- Inference:
 - marginalization, conditioning, completion
 - need for an efficient method (will discuss it)
- Learning: MLE if labeled examples are given

HMM POS Example

To understand better issues involved in efficient HMM inference, we will use a very small, walk-through example in POS tagging. We assume that the hidden internal states of the HMM correspond to correct POS tags of words, while the words correspond to generated observed variables. According to this, a sentence of n words would be associated with the following HMM graph in the unrolled form:



The variables W_1, \dots, W_n are assigned to words in the sentence, while variables T_1, \dots, T_n are assigned POS tags. The three probability tables that we mentioned in the definition of HMM are: $P(T_1)$, $P(T_{i+1}|T_i)$, and $P(W_i|T_i)$.

Slide notes:

Learning HMM (Training)

- Let us see how to learn HMM from a small set of these two labeled sentences:

```

swat V flies N like P ants N
time N flies V like P an D arrow N
  
```