

sequence modeling: part-of-speech tagging with hidden Markov models

CS 585, Fall 2019

Introduction to Natural Language Processing
<http://people.cs.umass.edu/~miyyer/cs585/>

Mohit Iyer

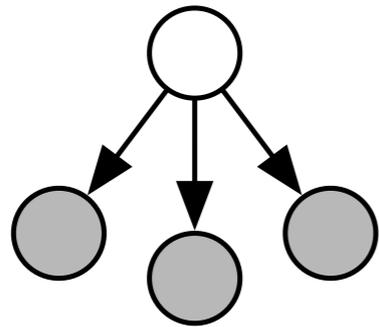
College of Information and Computer Sciences
University of Massachusetts Amherst

many slides from Brendan O'Connor & Jordan Boyd-Graber

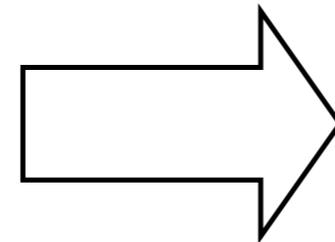
questions from last time...

- Busy next 2 weeks!
 - HW2! Due tmrw
 - Project milestone 1: due Oct 24
 - Midterm: Oct 31
- tested on optional readings? no
- final presentations? possibly Dec 12
- stats on HWs?
- what is a tensor?

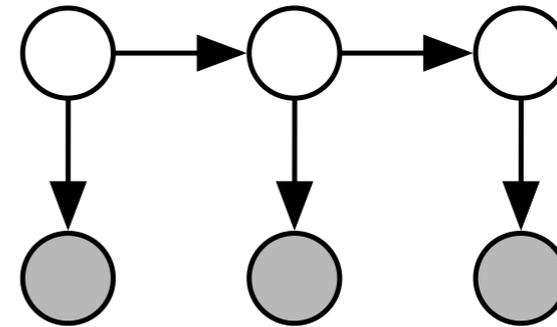
These are all **log-linear** models



Naive Bayes



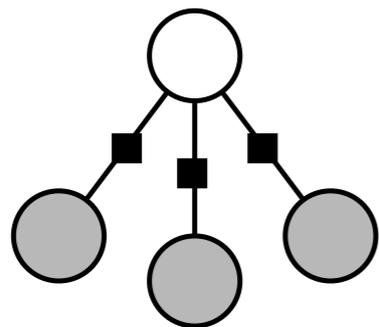
SEQUENCE



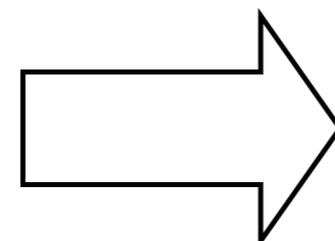
HMMs



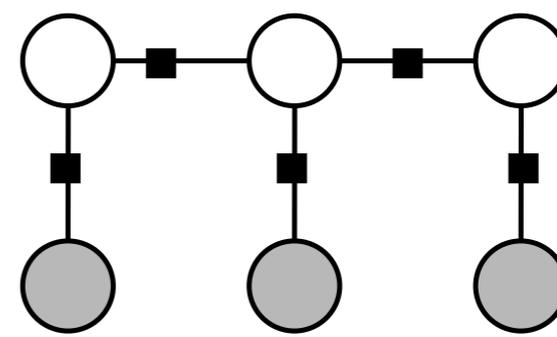
CONDITIONAL



Logistic Regression

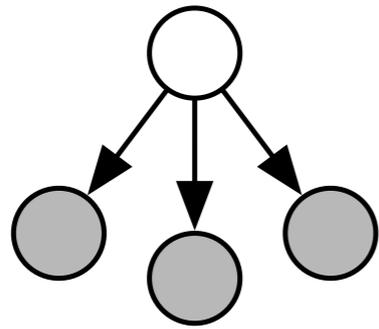


SEQUENCE

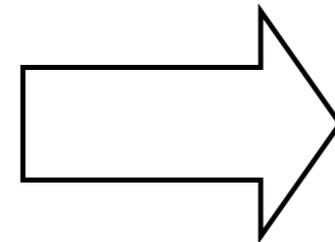


Linear-chain CRFs

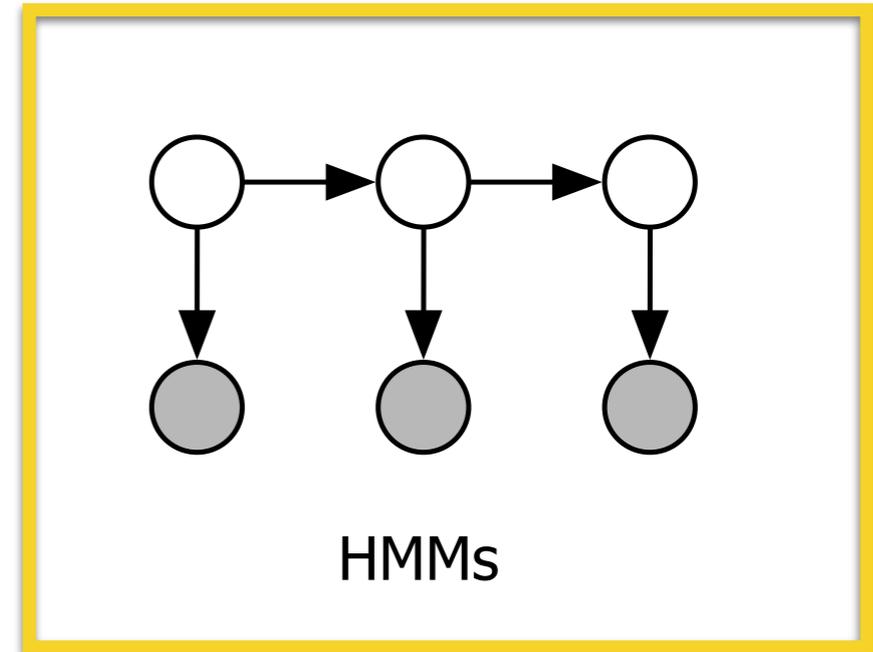
These are all **log-linear** models



Naive Bayes



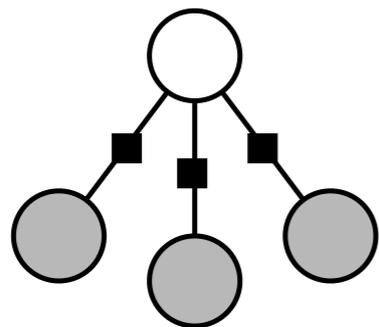
SEQUENCE



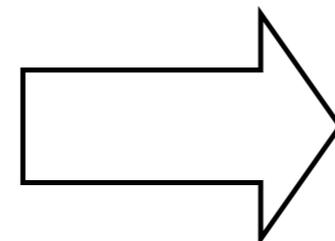
HMMs



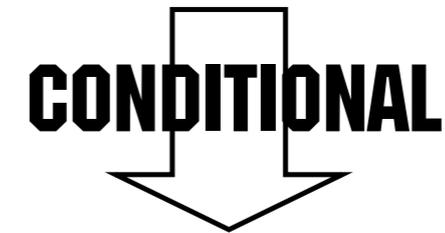
CONDITIONAL



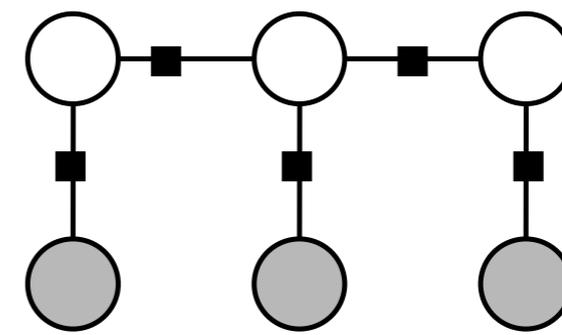
Logistic Regression



SEQUENCE

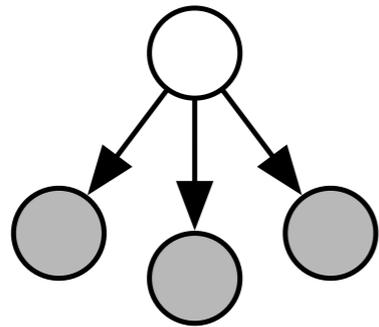


CONDITIONAL

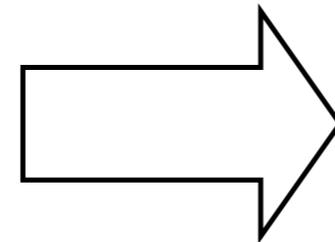


Linear-chain CRFs

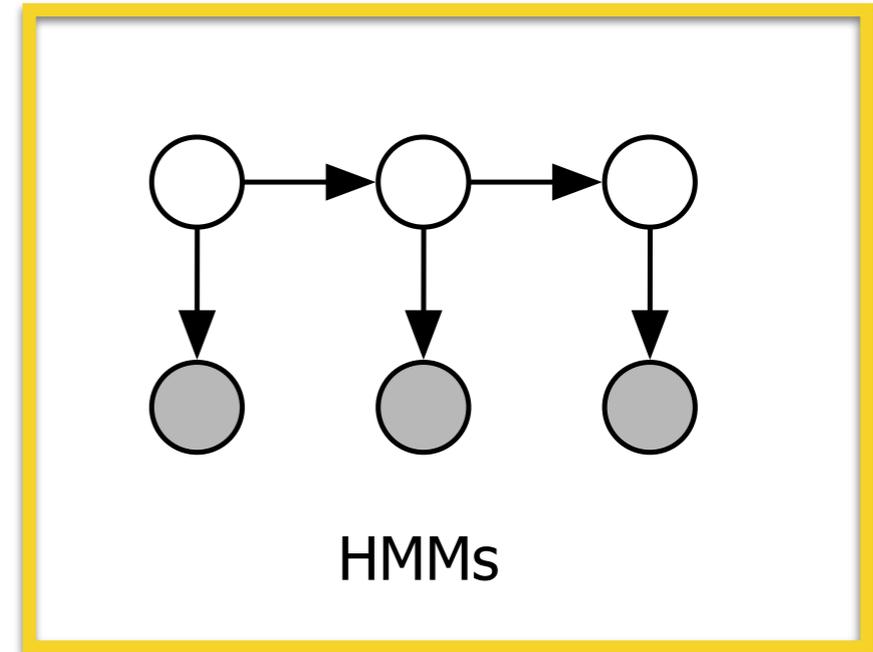
These are all **log-linear** models



Naive Bayes



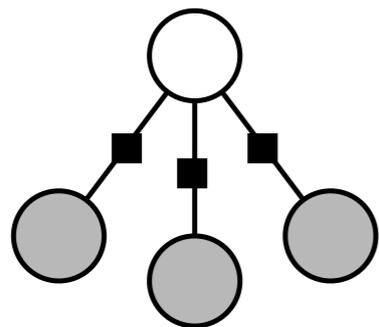
SEQUENCE



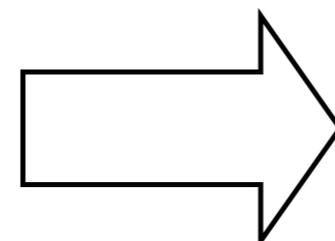
HMMs



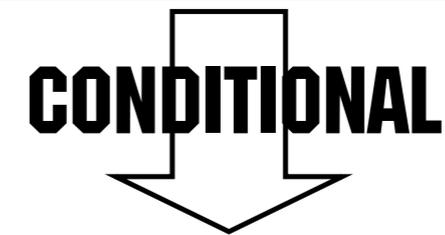
CONDITIONAL



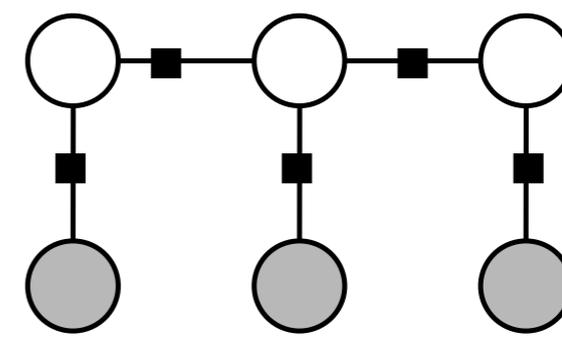
Logistic Regression



SEQUENCE



CONDITIONAL



Linear-chain CRFs

are neural networks log-linear models?

Tagging (Sequence Labeling)

- Given a sequence (in NLP, words), assign appropriate labels to each word.
- Many NLP problems can be viewed as sequence labeling:
 - POS Tagging
 - Chunking
 - Named Entity Tagging
- Labels of tokens are dependent on the labels of other tokens in the sequence, particularly their neighbors

Plays well with others.

VBZ RB IN NNS

What's a part-of-speech (POS)?

- Syntax = how words compose to form larger meaning bearing units
- POS = syntactic categories for words (a.k.a word class)
 - You could substitute words within a class and have a syntactically valid sentence

I saw the **dog**

I saw the **cat**

I saw the **___**

- Gives information how words combine into larger phrases

Why do we want POS?

- Useful for many syntactic and other NLP tasks.
 - Phrase identification (“chunking”)
 - Named entity recognition
 - Full parsing
 - Sentiment
- Especially when there’s a low amount of training data

POS patterns: sentiment

- Turney (2002): identify bigram phrases, from unlabeled corpus, useful for sentiment analysis.

Table 1. Patterns of tags for extracting two-word phrases from reviews.

| | First Word | Second Word | Third Word (Not Extracted) |
|----|--------------------|-------------------------|-------------------------------|
| 1. | JJ | NN or NNS | anything |
| 2. | RB, RBR, or RBS | JJ | not NN nor NNS |
| 3. | JJ | JJ | not NN nor NNS |
| 4. | NN or NNS | JJ | not NN nor NNS |
| 5. | RB, RBR, or RBS | VB, VBD, VBN, or VBG | anything |



Table 2. An example of the processing of a review that the author has classified as *recommended*.⁶

| Extracted Phrase | Part-of-Speech Tags | Semantic Orientation |
|-------------------|------------------------|-------------------------|
| online experience | JJ NN | 2.253 |
| low fees | JJ NNS | 0.333 |
| local branch | JJ NN | 0.421 |
| small part | JJ NN | 0.053 |
| online service | JJ NN | 2.780 |
| printable version | JJ NN | -0.705 |
| direct deposit | JJ NN | 1.288 |
| well other | RB JJ | 0.237 |
| inconveniently | RB VBN | -1.541 |
| located | | |
| other bank | JJ NN | -0.850 |
| true service | JJ NN | -0.732 |

POS patterns: simple noun phrases

- Quick and dirty noun phrase identification

| Tag Pattern | Example |
|-------------|---|
| A N | <i>linear function</i> |
| N N | <i>regression coefficients</i> |
| A A N | <i>Gaussian random variable</i> |
| A N N | <i>cumulative distribution function</i> |
| N A N | <i>mean squared error</i> |
| N N N | <i>class probability function</i> |
| N P N | <i>degrees of freedom</i> |

Table 5.2 Part of speech tag patterns for collocation filtering. These patterns were used by Justeson and Katz to identify likely collocations among frequently occurring word sequences.

Open class (lexical) words

Nouns

Proper

IBM
Italy

Common

cat / cats
snow

Verbs

Main

see
registered

Adjectives *old older oldest*

Adverbs *slowly*

Numbers

122,312
one

... more

Closed class (functional)

Determiners *the some*

Conjunctions *and or*

Pronouns *he its*

Modals

can
had

Prepositions *to with*

Particles *off up*

... more

Interjections *Ow Eh*

Open vs. Closed classes

- Open vs. Closed classes
 - Closed:
 - determiners: *a, an, the*
 - pronouns: *she, he, I*
 - prepositions: *on, under, over, near, by, ...*
 - Q: why called “closed”?
 - Open:
 - Nouns, Verbs, Adjectives, Adverbs.

Many Tagging Standards

- Penn Treebank (45 tags) ... this is the most common one
- Brown corpus (85 tags)
- Coarse tagsets
 - Universal POS tags (Petrov et. al. <https://github.com/slavpetrov/universal-pos-tags>)
 - Motivation: cross-linguistic regularities

Penn Treebank POS

- 45 possible tags
- 34 pages of tagging guidelines

| Tag | Description | Example | Tag | Description | Example |
|-------|-----------------------|------------------------|-----|-----------------------|---------------------------|
| CC | Coordin. Conjunction | <i>and, but, or</i> | SYM | Symbol | <i>+, %, &</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> | VCN | 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> | WPS | 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>([, (, { , <</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-final punc | <i>(. ! ?)</i> |
| RBS | Adverb, superlative | <i>fastest</i> | : | Mid-sentence punc | <i>(: ; ... - -)</i> |
| RP | Particle | <i>up, off</i> | | | |

<https://catalog ldc.upenn.edu/docs/LDC99T42/tagguid1.pdf>

Ambiguity in POS Tagging

- Words often have more than one POS: *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 POS tag for a particular instance of a word.

POS Tagging

- Input: Plays well with others
- Ambiguity: NNS/VBZ UH/JJ/NN/RB IN NNS
- Output: Plays/VBZ well/RB with/IN others/NNS

Penn
Treebank
POS tags

POS Tagging Performance

- How many tags are correct? (Tag Accuracy)
 - About 97% currently
 - But baseline is already 90%
 - Baseline is performance of stupidest possible method
 - Tag every word with its most frequent tag
 - Tag unknown words as nouns
- Partly easy because
 - Many words are unambiguous
 - You get points for them (*the*, *a*, etc.) and for punctuation marks!

How difficult is POS tagging?

- About 11% of the **word types** in the Brown corpus are ambiguous with regard to part of speech
- But they tend to be very common words. E.g., *that*
 - I know *that* he is honest = IN
 - Yes, *that* play was nice = DT
 - You can't go *that* far = RB
- 40% of the **word tokens** are ambiguous

Token vs. Type

Token is instance or individual occurrence of a type.

Two Types of Constraints

Influential/JJ members/NNS of/IN the/DT House/NNP Ways/NNP and/CC Means/NNP Committee/NNP introduced/VBD legislation/NN that/WDT would/MD restrict/VB how/WRB the/DT new/JJ savings-and-loan/NN bailout/NN agency/NN can/MD raise/VB capital/NN ./.

- ▶ “Local”: e.g., *can* is more likely to be a modal verb MD rather than a noun NN
- ▶ “Contextual”: e.g., a noun is much more likely than a verb to follow a determiner
- ▶ Sometimes these preferences are in conflict:

The trash can is in the garage

Generative Model

- Probabilistic generative model for sequences.
- Assume an underlying set of hidden (unobserved) states in which the model can be (e.g. parts of speech). **different from RNN hidden states!**
- Assume probabilistic transitions between states over time (e.g. transition from POS to another POS as sequence is generated).
- Assume a probabilistic generation of tokens from states (e.g. words generated for each POS).

Hidden Markov Models

- ▶ We have an input sentence $x = x_1, x_2, \dots, x_n$
(x_i is the i 'th word in the sentence)
- ▶ We have a tag sequence $y = y_1, y_2, \dots, y_n$
(y_i is the i 'th tag in the sentence)

- ▶ We'll use an HMM to define

$$p(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n)$$

for any sentence $x_1 \dots x_n$ and tag sequence $y_1 \dots y_n$ of the same length.

- ▶ Then the most likely tag sequence for x is

$$\arg \max_{y_1 \dots y_n} p(x_1 \dots x_n, y_1, y_2, \dots, y_n)$$

are HMMs generative or discriminative models?

HMM Definition

Assume K parts of speech, a lexicon size of V , a series of observations $\{x_1, \dots, x_N\}$, and a series of unobserved states $\{z_1, \dots, z_N\}$.

π A distribution over start states (vector of length K):

$$\pi_i = p(z_1 = i)$$

θ Transition matrix (matrix of size K by K):

$$\theta_{i,j} = p(z_n = j | z_{n-1} = i) \quad \text{Markov assumption!}$$

β An emission matrix (matrix of size K by V):

$$\beta_{j,w} = p(x_n = w | z_n = j)$$

HMM Definition

Assume K parts of speech, a lexicon size of V , a series of observations $\{x_1, \dots, x_N\}$, and a series of unobserved states $\{z_1, \dots, z_N\}$.

π A distribution over start states (vector of length K):

$$\pi_i = p(z_1 = i)$$

θ Transition matrix (matrix of size K by K):

$$\theta_{i,j} = p(z_n = j | z_{n-1} = i)$$

β An emission matrix (matrix of size K by V):

$$\beta_{j,w} = p(x_n = w | z_n = j)$$

Two problems: How do we move from data to a model? (Estimation)

How do we move from a model and unlabeled data to labeled data?

(Inference)

today: estimation

Reminder: How do we estimate a probability?

- For a multinomial distribution (i.e. a discrete distribution, like over words):

$$\theta_i = \frac{n_i + \alpha_i}{\sum_k n_k + \alpha_k} \quad (1)$$

- α_i is called a smoothing factor, a pseudocount, etc.

just like in naive Bayes, we'll be counting to estimate these probabilities!

Training Sentences

x = tokens

z = POS tags

x here come old flattop
z MOD V MOD N

a crowd of people stopped and stared
DET N PREP N V CONJ V

gotta get you into my life
V V PRO PREP PRO V

and I love her
CONJ PRO V PRO

Initial Probability π

| POS | Frequency | Probability |
|------|-----------|-------------|
| MOD | 1.1 | 0.234 |
| DET | 1.1 | 0.234 |
| CONJ | 1.1 | 0.234 |
| N | 0.1 | 0.021 |
| PREP | 0.1 | 0.021 |
| PRO | 0.1 | 0.021 |
| V | 1.1 | 0.234 |

let's use add-alpha smoothing with $\alpha = 0.1$

Training Sentences

here come old flattop
MOD V MOD N

a crowd of people stopped and stared
DET N PREP N V CONJ V

gotta get you into my life
V V PRO PREP PRO N

and I love her
CONJ PRO V PRO

Training Sentences

here come old flattop
MOD V MOD N

a crowd of people stopped and stared
DET N PREP N V CONJ V

gotta get you into my life
V V PRO PREP PRO N

and I love her
CONJ PRO V PRO

Training Sentences

here come old flattop
MOD V MOD N

a crowd of people stopped and stared
DET N PREP N V CONJ V

gotta get you into my life
V V PRO PREP PRO N

and I love her
CONJ PRO V PRO

Transition Probability θ

- We can ignore the words; just look at the parts of speech. Let's compute one row, the row for verbs.
- We see the following transitions: $V \rightarrow \text{MOD}$, $V \rightarrow \text{CONJ}$, $V \rightarrow V$, $V \rightarrow \text{PRO}$, and $V \rightarrow \text{PRO}$

| POS | Frequency | Probability |
|------|-----------|-------------|
| MOD | 1.1 | 0.193 |
| DET | 0.1 | 0.018 |
| CONJ | 1.1 | 0.193 |
| N | 0.1 | 0.018 |
| PREP | 0.1 | 0.018 |
| PRO | 2.1 | 0.368 |
| V | 1.1 | 0.193 |

how many transition probability distributions do we have?

Training Sentences

here come old flattop
MOD V MOD N

a crowd of people stopped and stared
DET N PREP N V CONJ V

gotta get you into my life
V V PRO PREP PRO N

and I love her
CONJ PRO V PRO

Training Sentences

here **come** old flattop
MOD V MOD N

a crowd of people **stopped** and **stared**
DET N PREP N V CONJ V

gotta **get** you into my life
V V PRO PREP PRO N

and I **love** her
CONJ PRO V PRO

Emission Probability β

Let's look at verbs ...

| | | | | | |
|-------------|--------|--------|--------|--------|---------|
| Word | a | and | come | crowd | flattop |
| Frequency | 0.1 | 0.1 | 1.1 | 0.1 | 0.1 |
| Probability | 0.0125 | 0.0125 | 0.1375 | 0.0125 | 0.0125 |
| Word | get | gotta | her | here | i |
| Frequency | 1.1 | 1.1 | 0.1 | 0.1 | 0.1 |
| Probability | 0.1375 | 0.1375 | 0.0125 | 0.0125 | 0.0125 |
| Word | into | it | life | love | my |
| Frequency | 0.1 | 0.1 | 0.1 | 1.1 | 0.1 |
| Probability | 0.0125 | 0.0125 | 0.0125 | 0.1375 | 0.0125 |
| Word | of | old | people | stared | stopped |
| Frequency | 0.1 | 0.1 | 0.1 | 1.1 | 1.1 |
| Probability | 0.0125 | 0.0125 | 0.0125 | 0.1375 | 0.1375 |

how many emission probability distributions do we have?

Next time ...

- Viterbi algorithm: dynamic algorithm discovering the most likely POS sequence given a sentence

what if we don't have any labeled data to estimate an HMM?
we can still learn a model using the expectation-maximization
algorithm. but we won't cover this in class :(