

Lecture 4: Viterbi, NER



Cornell Bowers CIS
Computer Science

Claire Cardie, Tanya Goyal

CS 4740 (and crosslists): Introduction to Natural Language Processing

Announcements

- ▶ HW1 released today! Due Fri Feb 21 11:59pm.
 - ▶ START IT NOW!!!!!!

Today

- HMMs as a tagging technology: Viterbi
 - You will implement for HW1!!!
- HMMs as a generative model
- Where do the probabilities come from?
- Named entity tagging: the task for HW1!!!

Recall: HMM POS Tagger

? ? ?
Cornell beat Harvard

Goal: Find the tag sequence that maximizes $P(t_1 \dots t_N \mid w_1 \dots w_N)$

Need to Bayes flip:

$$= \frac{P(w_1 \dots w_N \mid t_1 \dots t_N) \cdot P(t_1 \dots t_N)}{\cancel{P(w_1 \dots w_N)}}$$

Make Independence and Markov Assumptions

? ? ?
Cornell beat Harvard

$$\underline{P(w_1 \dots w_N | t_1 \dots t_N) \cdot P(t_1 \dots t_N)}$$

$P(t_1, \dots, t_n)$: approximate using **n-gram model**

bigram $\prod_{i=1, n} P(t_i | t_{i-1})$

trigram $\prod_{i=1, n} P(t_i | t_{i-2} t_{i-1})$

Make Independence and Markov Assumptions

? ? ?
Cornell beat Harvard

$$\underline{P(w_1 \dots w_N | t_1 \dots t_N) P(t_1 \dots t_N)}$$

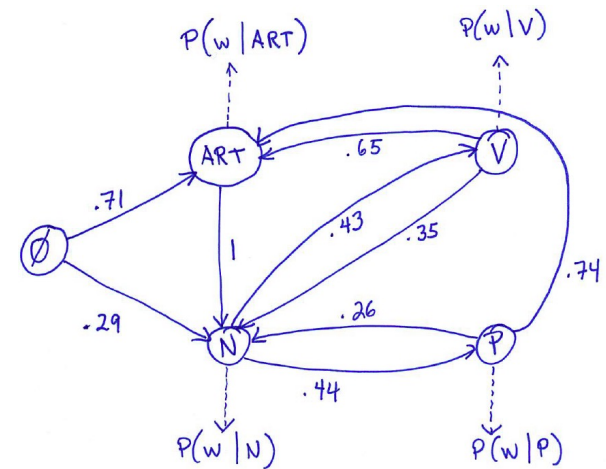
Assume each word appears with a particular tag independent of its neighbors

$$P(w_1 \dots w_n | t_1 \dots t_n) \cong \prod_{i=1,n} P(w_i | t_i)$$

? ? ?
Cornell beat Harvard

$$P(t_1 \dots t_N \mid w_1 \dots w_N) \cong \prod_{i=1, n} P(t_i \mid t_{i-1}) \cdot P(w_i \mid t_i)$$

- ▶ Equation is modeled by an HMM (probabilistic finite-state machine)
 - ▶ **States:** represent the possible POS
 - ▶ **Transition probabilities:** bigram probabilities for tags
 - ▶ **Emission (observation) probabilities:** indicate, for each word, how likely that word is to be selected if we randomly select a POS



Tagging algorithm

N V N
Cornell beat Harvard

Given a new sentence to tag

- For every possible tag sequence,
 - Apply equation to calculate the score
- Select the highest-scoring tag sequence

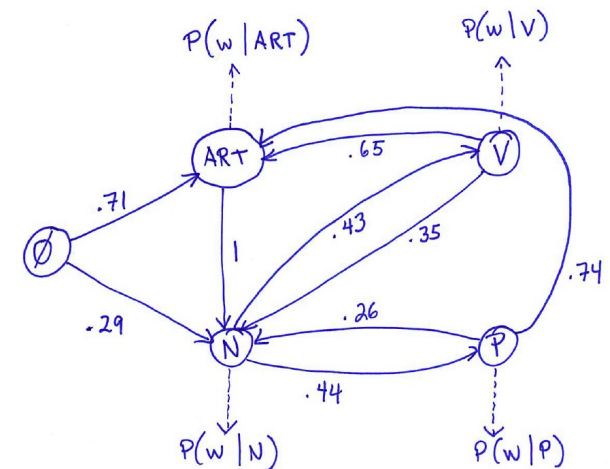
$$\prod_{i=1,n} P(t_i | t_{i-1}) \cdot P(w_i | t_i)$$

Uh-oh...Too many possible tag sequences to do this!!!

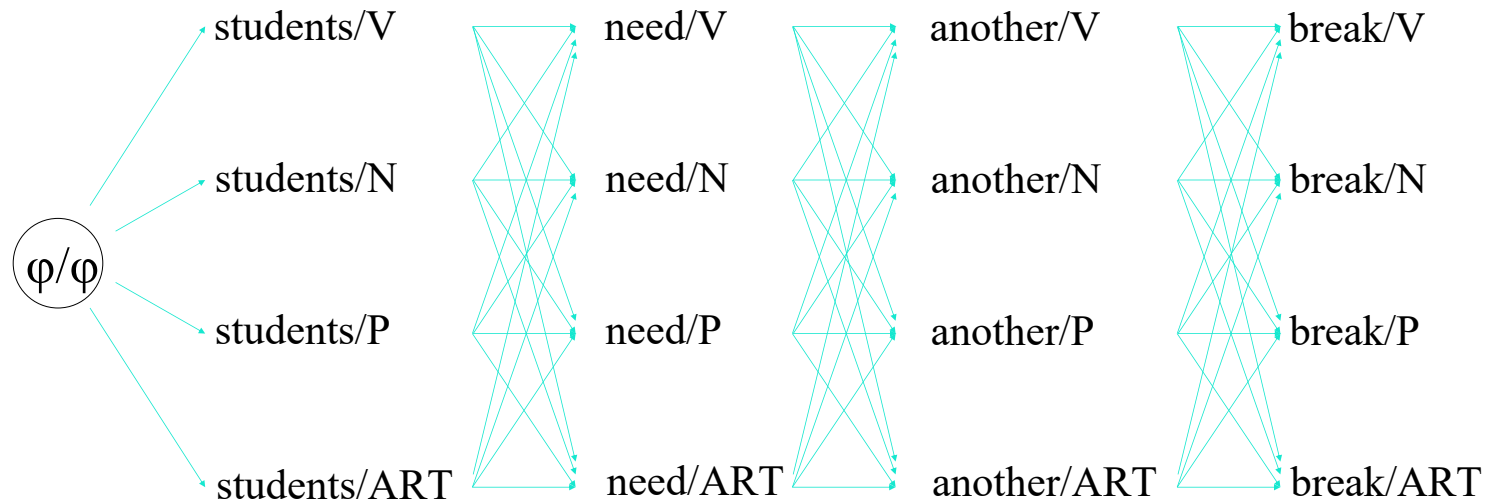
Sentence length $m=20$

Tagset of size $T = 15$

$T^m = 15^{20}$ tag sequences!!!

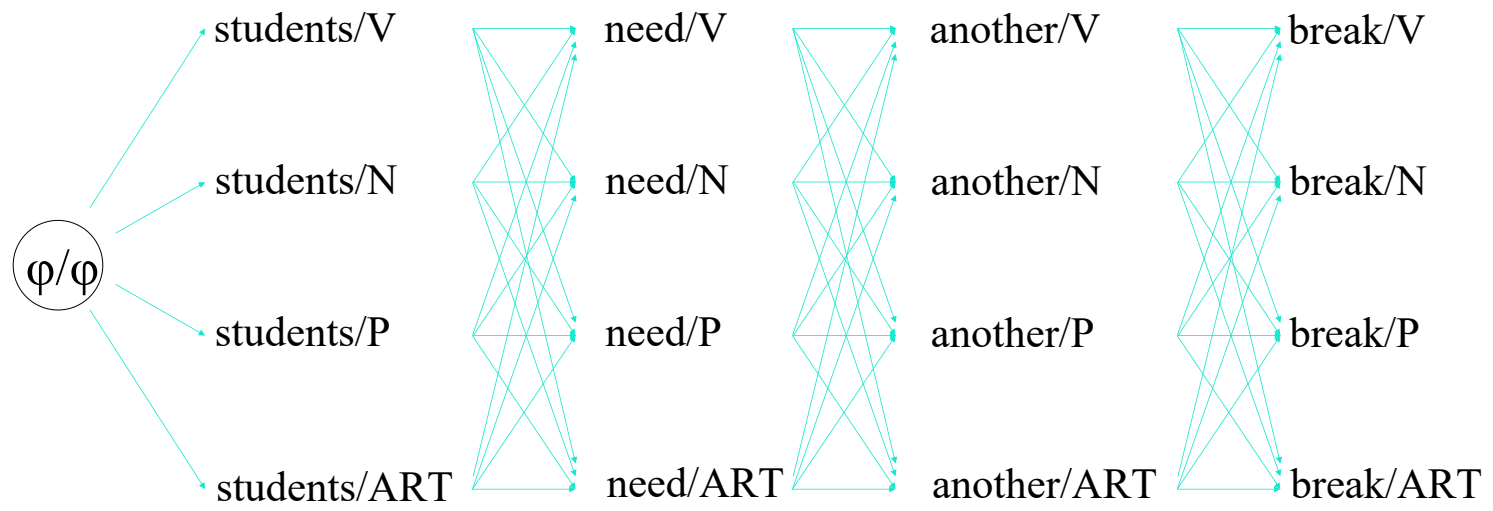


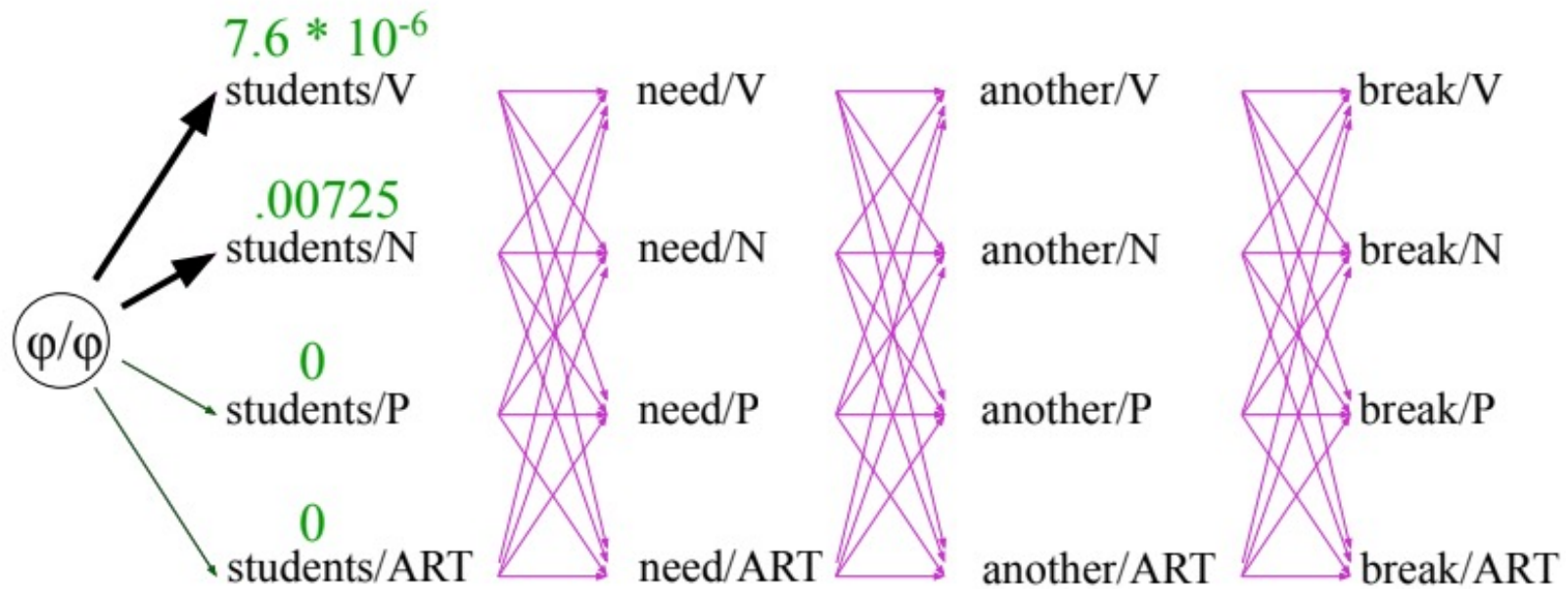
How do we avoid computing the probabilities for all possible paths?



Viterbi Algorithm Allows Efficient Search for the Most Likely Sequence

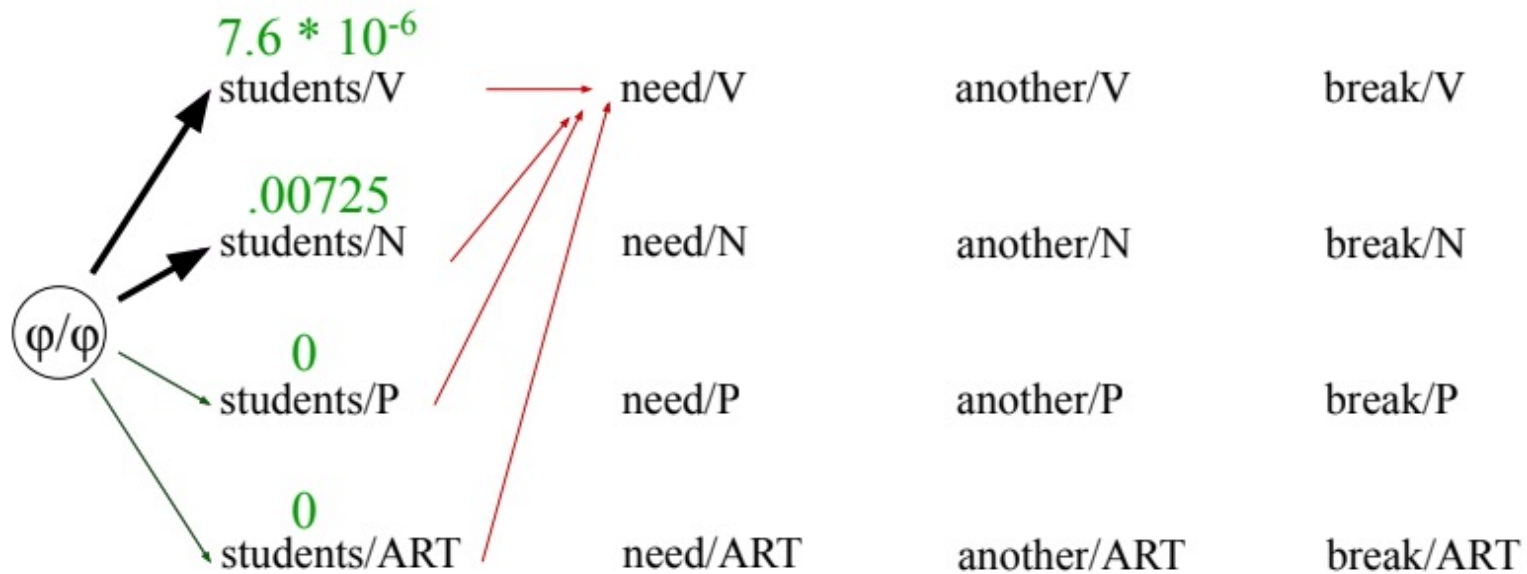
- Key idea: Markov assumptions mean that we do not need to enumerate all possible sequences
- Viterbi algorithm
 - Sweep forward, one word at a time, finding the most likely (highest-scoring) tag sequence ending with each possible tag
 - With the right bookkeeping, we can then “read off” the most likely tag sequence once we reach the end of the sentence





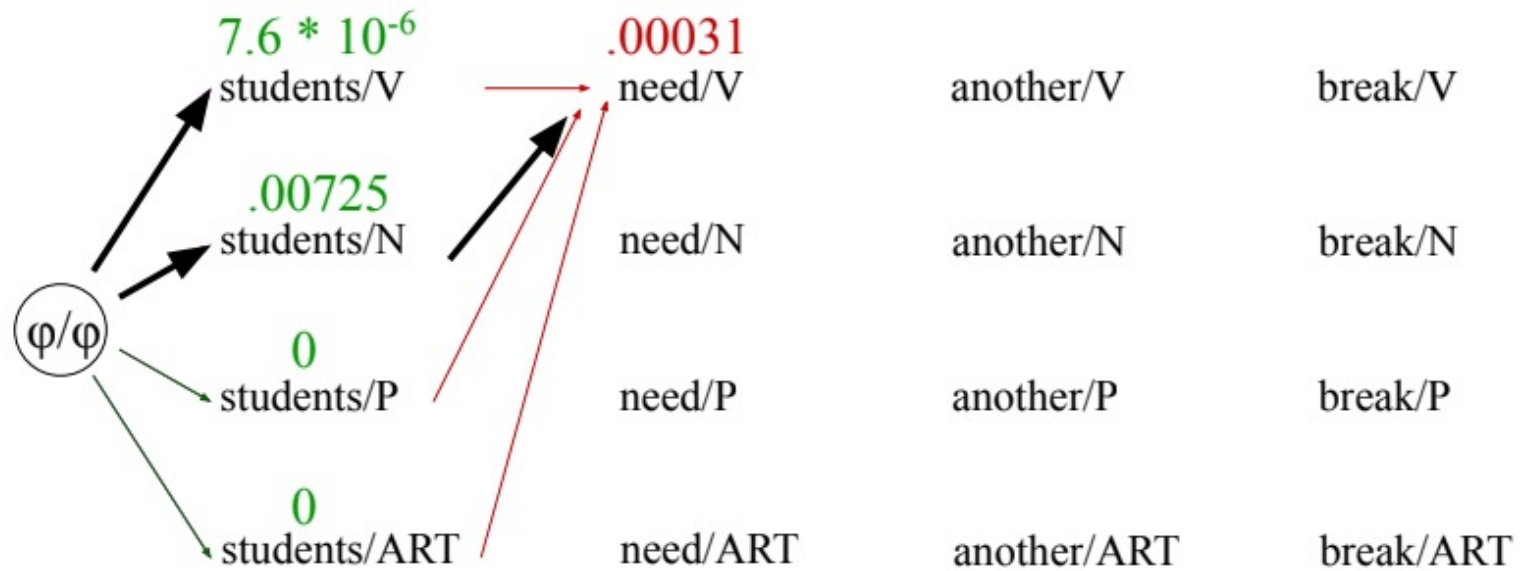
$$v_0(j) = P(\text{tag}_0=j \mid \langle s \rangle) * P(\text{students} \mid \text{tag}_0=j)$$

$$vb_0(j) = \langle s \rangle$$



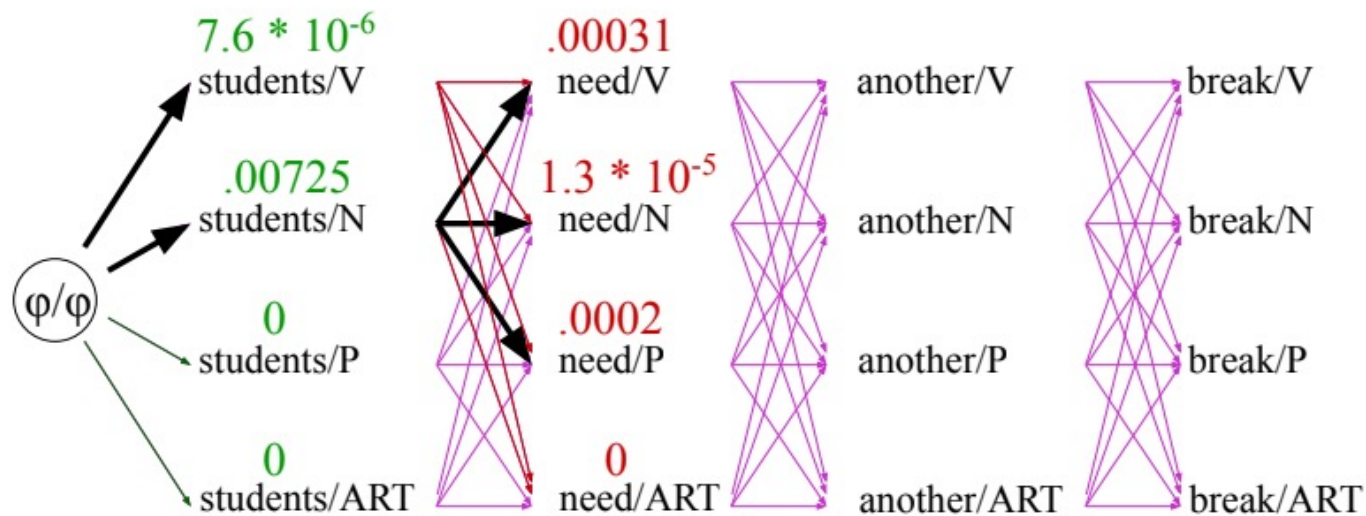
$$v_1(j) = \max_{i=1}^J v_0(i) * P(\text{tag}_1=j \mid \text{tag}_0=i) * P(\text{need} \mid \text{tag}_1=j)$$

$$vb_t(j) = \text{prev tag that maximizes } v_t(j)$$



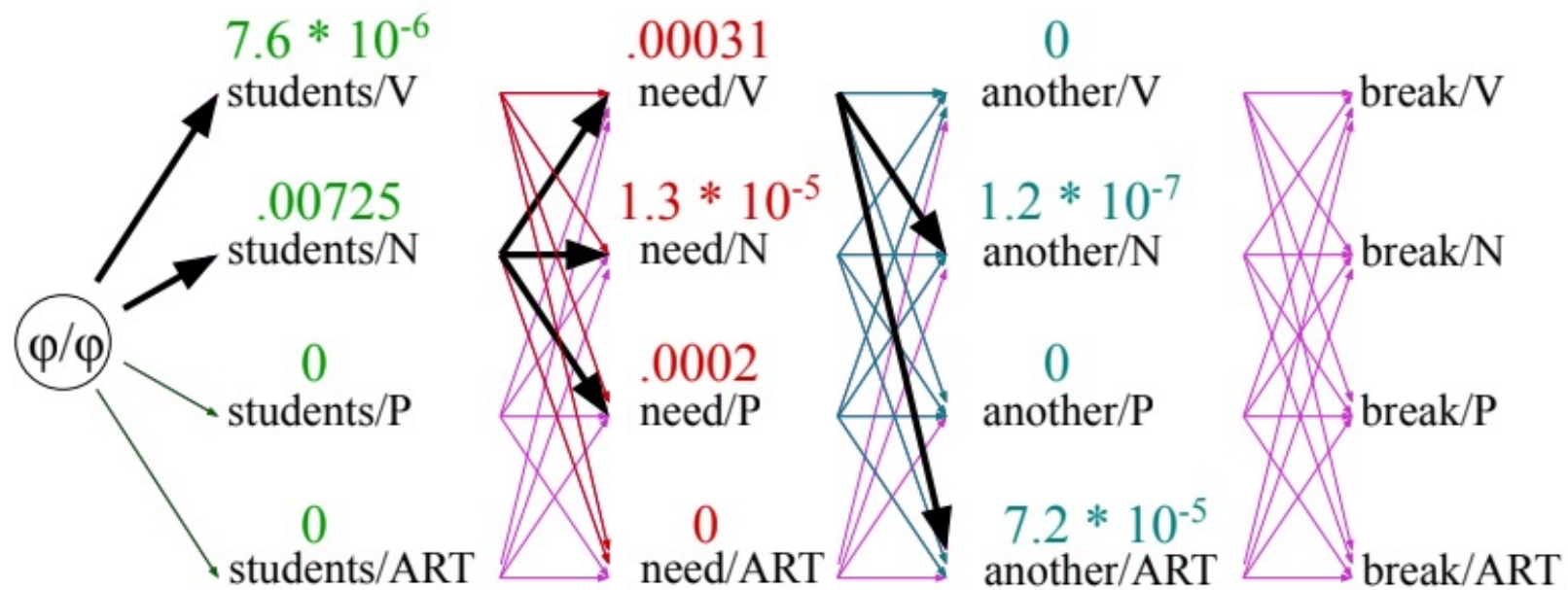
$$v_1(j) = \max_{i=1}^J v_0(i) * P(\text{tag}_1=j \mid \text{tag}_0=i) * P(\text{need} \mid \text{tag}_1=j)$$

$$vb_t(j) = \text{prev tag that maximizes } v_t(j)$$



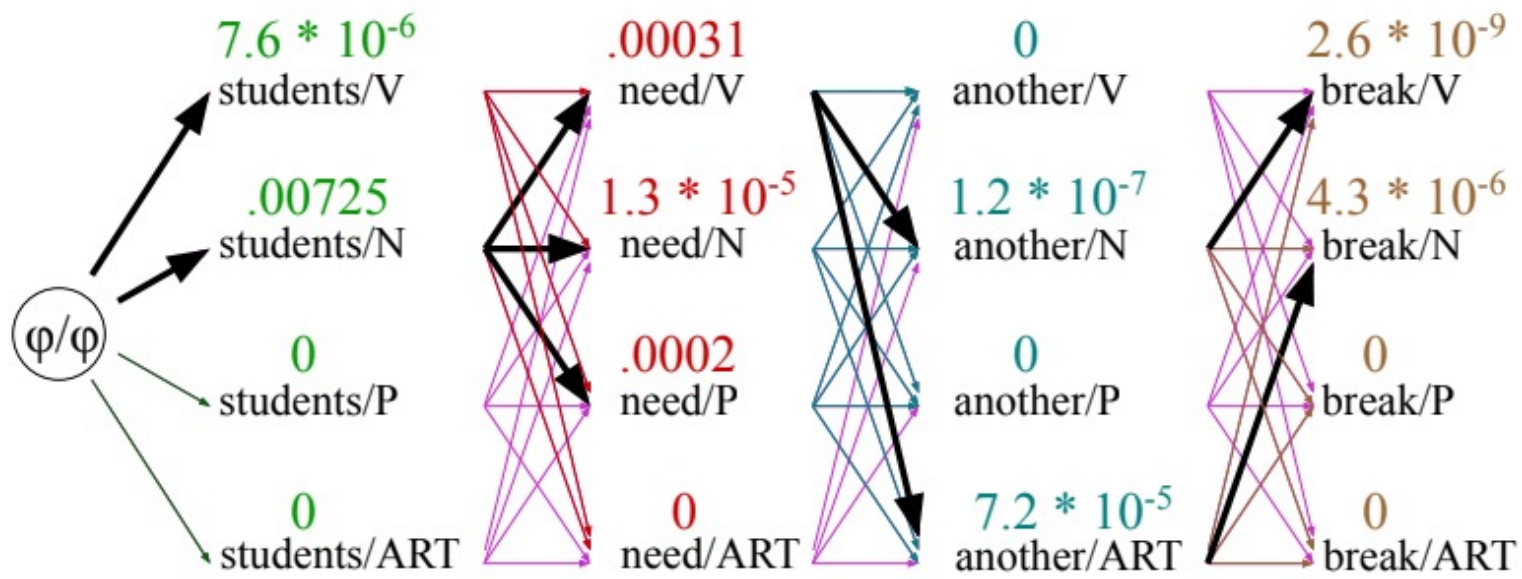
$$v_1(j) = \max_{i=1}^J v_0(i) * P(\text{tag}_1=j \mid \text{tag}_0=i) * P(\text{need} \mid \text{tag}_1=j)$$

$$vb_t(j) = \text{prev tag that maximizes } v_t(j)$$



$$v_2(j) = \max_{i=1}^J v_1(i) * P(\text{tag}_2=j \mid \text{tag}_1=i) * P(\text{another} \mid \text{tag}_2=j)$$

$$vb_t(j) = \text{prev tag that maximizes } v_t(j)$$

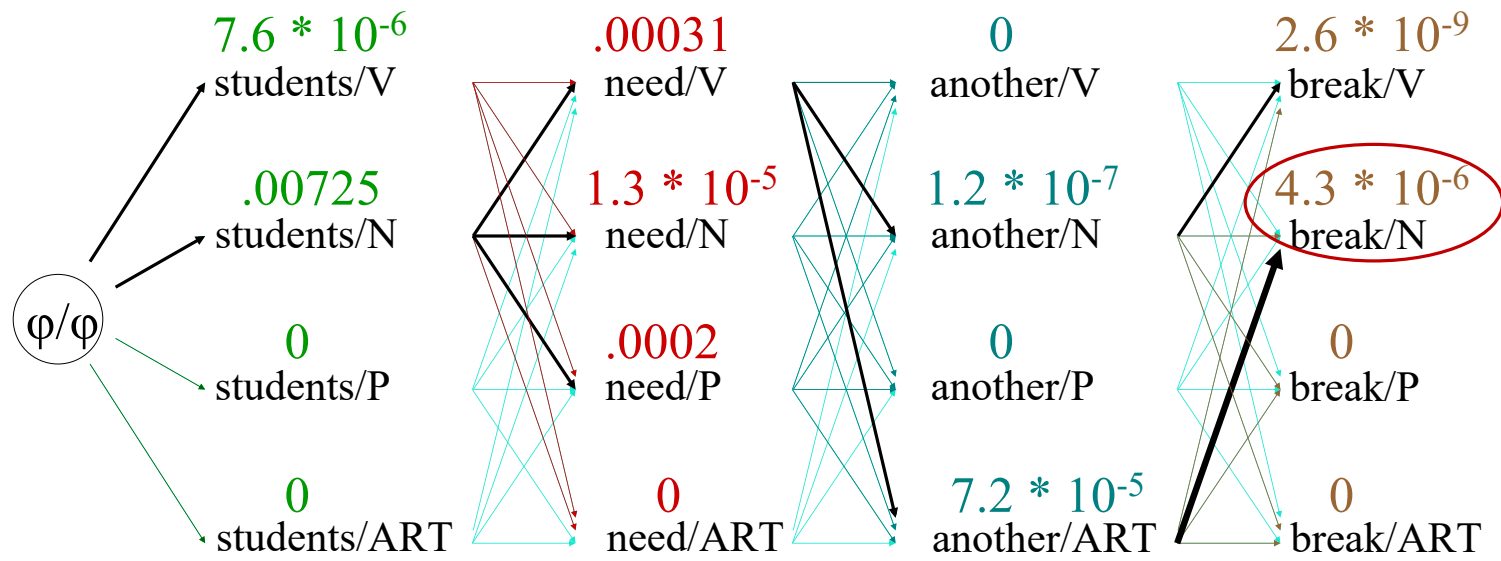


$$v_3(j) = \max_{i=1}^J v_2(i) * P(\text{tag}_3=j \mid \text{tag}_2=i) * P(\text{break} \mid \text{tag}_3=j)$$

$$vb_t(j) = \text{prev tag that maximizes } v_t(j)$$

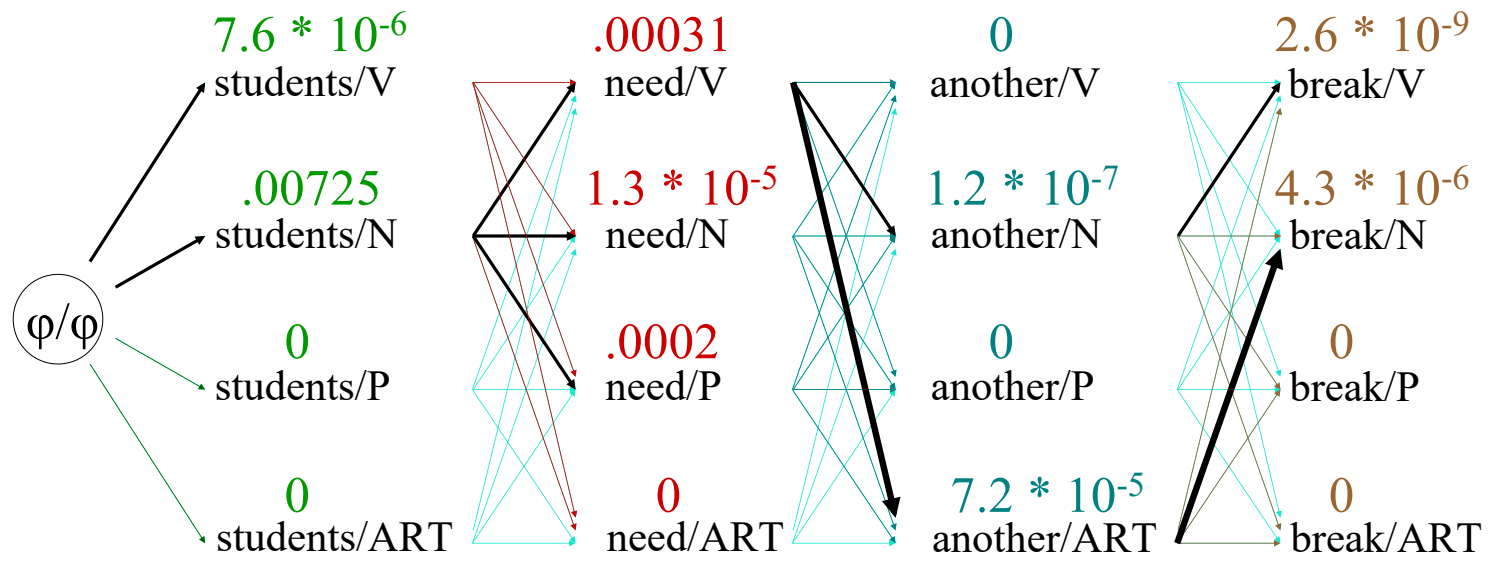
- To assign the maximum probability tag sequence, follow the backpointers that led to the largest product at v_3 !

Viterbi Algorithm



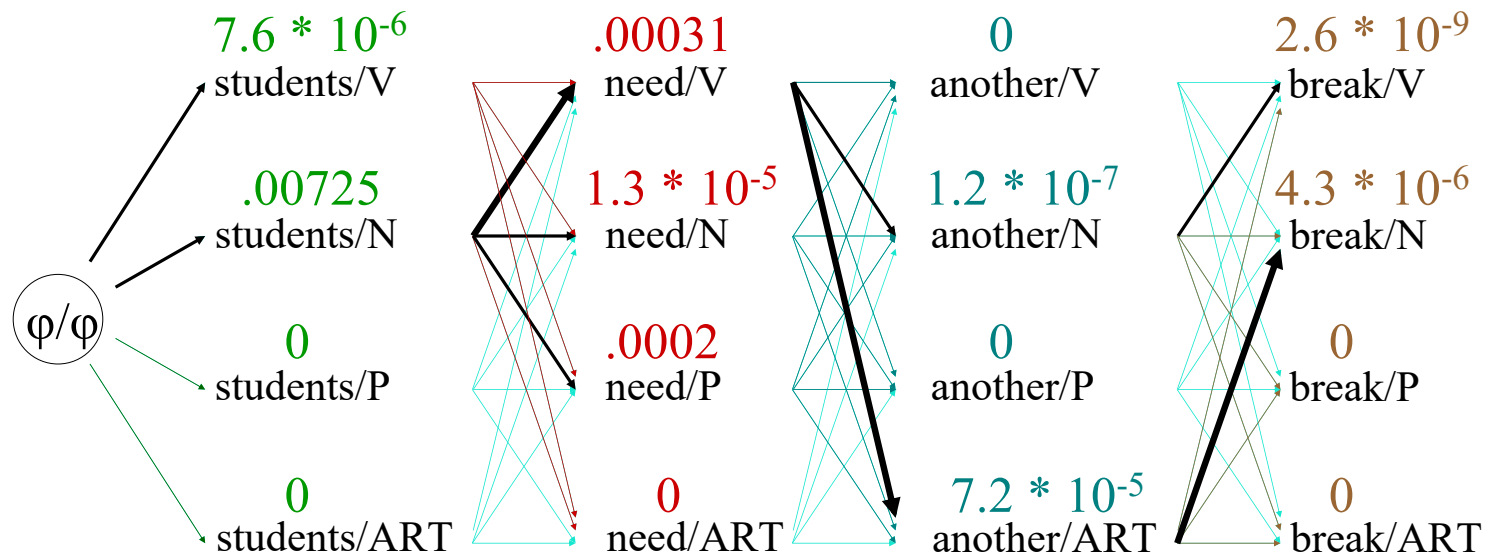
$t_3 = N$

Viterbi Algorithm



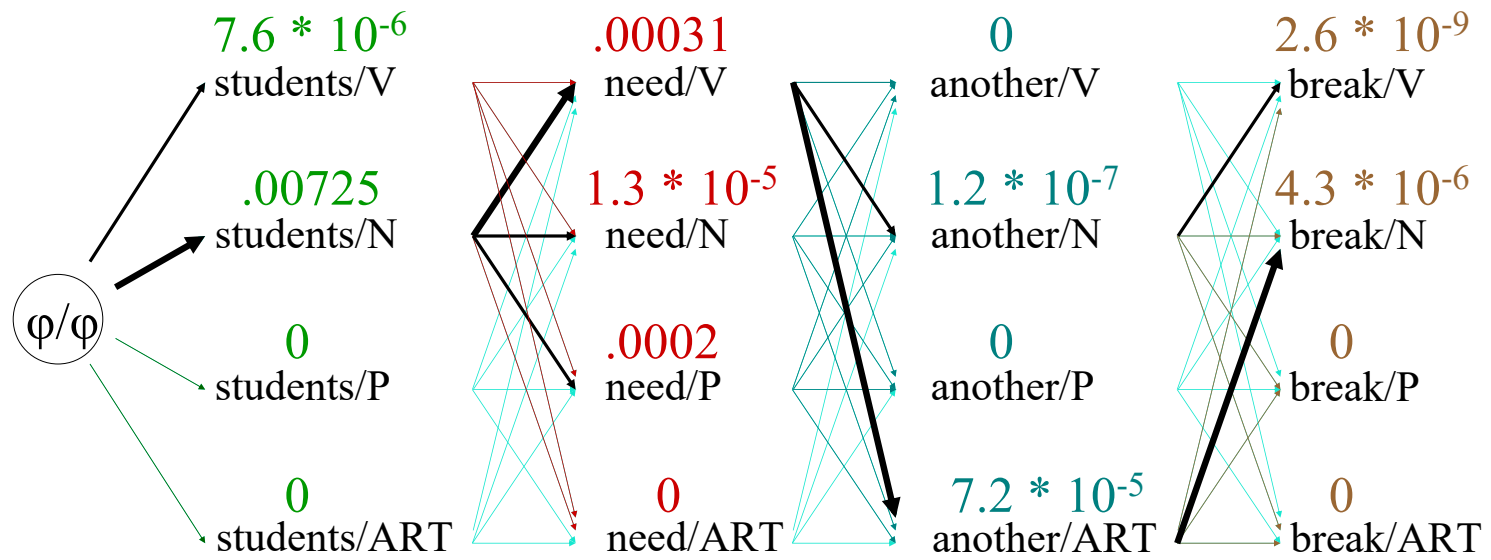
$t_3 = N, t_4 = ART$

Viterbi Algorithm



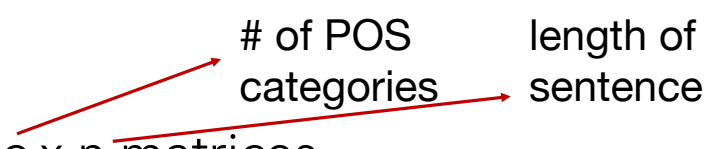
t₃= N, t₂= ART, t₁= V

Viterbi Algorithm



$t_3 = N, t_2 = ART, t_1 = V, t_0 = N$

Time/space complexity

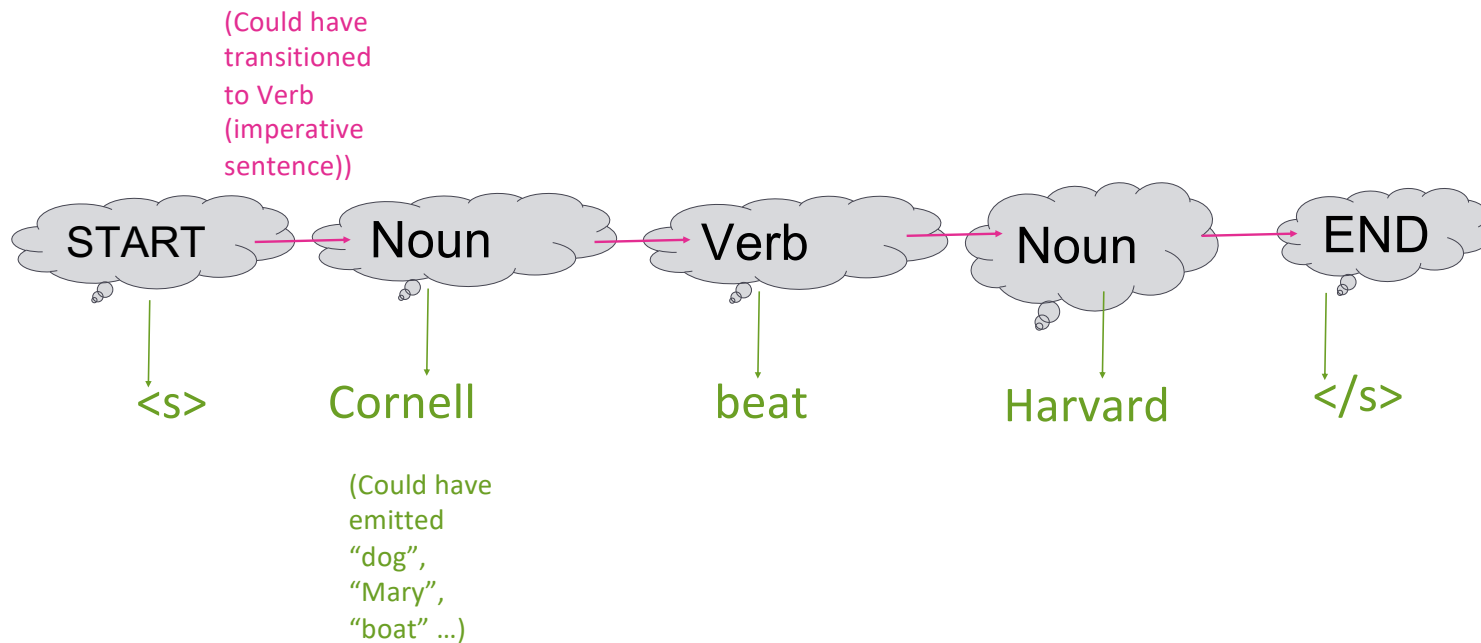
- Space
 - Two $c \times n$ matrices
 - (and data structure for transition and lexical generation probabilities)
 - Time
 - $O(c^2n)$ for forward pass
 - $O(n)$ for backward pass
 - Much better than the $O(c^n)$ brute force option
- 

Today

- HMMs as a tagging technology: Viterbi
 - You will implement for HW1!!!
- **HMMs as a generative model**
- Where do the probabilities come from?
- Named entity tagging: the task for HW1!!!

HMMs as sentence generators

When in an underlying state (POS), generate a token.
Then, choose a next underlying state.



Today

- HMMs as a tagging technology: Viterbi
 - You will implement for HW1!!!
- HMMs as a generative model
- **Where do the probabilities come from?**
- Named entity tagging: the task for HW1!!!

Where do HMM transitions/emission probs come from?

Assume that we have *labelled data*:

For every observed token x_i , the (usually hidden) true tag c_i is given.

<s>/<s> I/PP am/VBP sitting/VBG in/IN Mindy/NNP 's/POS restaurant/NN
eating/VBG the/DT gefilte/NN fish/NN ./</s>/</s>

- Looks like VBG generates things like “sitting” and “eating”; and a period (.) can be followed by </s>.

Warning: training data might omit <s>, <s>, </s>, </s>. You'll want to insert them (implicitly or explicitly).

“Raw count” method for setting transition and emission probs

$$P_{HMM}(w_j | c) := \frac{\text{count (word } w_j \text{ in training with tag } c)}{\text{count (word tokens in training with tag } c)}$$

$$P_{HMM}(c' | c) := \frac{\text{count (} c \text{ followed by } c')}{\text{count (} c)}$$

Smoothing: “lack of evidence is not evidence of lack”

An unseen event isn't necessarily impossible! Safer to have all probs be non-zero.

One common smoothing technique: add-k.

$$P(\mathbf{b} \mid \mathbf{a}) := [\text{Count}(\mathbf{a} \mathbf{b}) + k] \dots$$

...divided by the normalization term:

$$\text{sum over all possible } \mathbf{b}' \text{ of } [C(\mathbf{a} \mathbf{b}') + k]$$

Today

- HMMs as a tagging technology: Viterbi
 - You will implement for HW1!!!
- HMMs as a generative model
- Where do the probabilities come from?
- **Named entity tagging: the task for HW1!!!**

Named Entity Recognition

Identify all:

- Named locations, named persons, named organizations, dates, times, monetary amounts...
- Fixed set of NE types

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	The Mt. Sanitas loop is in Sunshine Canyon .
Geo-Political Entity	GPE	countries, states, provinces	Palo Alto is raising the fees for parking.
Facility	FAC	bridges, buildings, airports	Consider the Golden Gate Bridge .
Vehicles	VEH	planes, trains, automobiles	It was a classic Ford Falcon .

Figure 17.1 A list of generic named entity types with the kinds of entities they refer to.

NER output

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

Named Entity Recognition (NER): note the multi-word named entities, like “North America”

In fact, the Chinese market has the three most influential names of the retail and tech space – Alibaba, Baidu, and Tencent (collectively touted as BAT), and is betting big in the global AI in retail industry space. The three giants which are claimed to have a cut-throat competition with the U.S. (in terms of resources and capital) are positioning themselves to become the ‘future AI platforms’. The trio is also expanding in other Asian countries and investing heavily in the U.S. based AI startups to leverage the power of AI. Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one, with an anticipated CAGR of 45% over 2018 - 2024.

Lots of errors!!!

To further elaborate on the geographical trends, North America has procured more than 50% of the global share in 2017 and has been leading the regional landscape of AI in the retail market. The U.S. has a significant credit in the regional trends with over 65% of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as Google, IBM, and Microsoft.

Ambiguity in NER

Name	Possible Categories
<i>Washington</i>	Person, Location, Political Entity, Organization, Vehicle
<i>Downing St.</i>	Location, Organization
<i>IRA</i>	Person, Organization, Monetary Instrument
<i>Louis Vuitton</i>	Person, Organization, Commercial Product

Figure 17.2 Common categorical ambiguities associated with various proper names.

[PER Washington] was born into slavery on the farm of James Burroughs.
[ORG Washington] went up 2 games to 1 in the four-game series.
Blair arrived in [LOC Washington] for what may well be his last state visit.
In June, [GPE Washington] passed a primary seatbelt law.
The [VEH Washington] had proved to be a leaky ship, every passage I made...

Figure 17.3 Examples of type ambiguities in the use of the name *Washington*.

NE Recognition

- Identify the text **spans** that correspond to the proper names (or dates, times, money expressions)
 - How do we describe a **chunk** of text using individual-word tags?
- Assign the correct named entity (NE) type

BIO tag set for NER

- Allows distinguishing adjacent NEs
 - We'll fly to **New Orleans** **Friday**
- B_{xxx} : First (ie. Beginning) token in an NE of type XXX
- I_{xxx} : Inside of an entity type XXX
- O : Outside of all NEs

BIO Tagging

B: token that *begins* a span

I: tokens *inside* a span

O: tokens outside of any span

of tags (where n is #entity types):

1 O tag,

n B tags,

n I tags

total of $2n+1$

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
.	O

HMMs for NE detection

Just like in POS tagging

- States **Q**
 - BIO tags
- Observations **O**
 - Word tokens
- Transition Probabilities **A**
 - $P(\text{BIOtag}_i | \text{BIOtag}_{i-1})$
- Emission (lexical generation) Probabilities **B**
 - $P(w_i | \text{BIOtag}_i)$

Find most likely BIO tag sequence using Viterbi

Reconstruct the NEs from the BIO tags

Take-aways

- HMMs as a tagging technology: the Viterbi algorithm for efficiently assigning the highest probability tag sequence
- HMMs as a generative model (just 1 slide)
- Where do the probabilities come from? Labeled data
- Named entity tagging: the task for HW1!!!