

CS4248: Natural Language Processing

Lecture 7 — Sequences

Announcements

Deadlines: all 14 Mar, 23:59

1. Assignment 2

- Goal: practice manual feature engineering.
- To be fair, only certain technologies already covered are allowed.
- 2. Midterm Survey 1% of your course marks (5 mins of your time)
- 3. Project: TEAMMATES Intra-Project Formative Feedback
 - (ungraded, but monitored)

4. Correction: Cross Entropy (Slide 34, Lecture 5): -clean version updated

Recap of Week 06



Pre-Lecture Activity for Next Week

- Assigned Task (due before Mar 8)
 - Post a 1-2 paragraph answer to the following question in the [Pre-Lecture] discussion
 - Watch the panel discussion of the recent Generative AI, Free Speech, & Public Discourse https://www.youtube.com/live/BBhewsinQwU?si=ODkIpYjqOCLZh8xD&t=8659



Relate an opinion by one of the panelists to the lecture materials presented today. Why did you pick this opinion to highlight and what is your own opinion on it?

Side notes:

- We will talk about this in the next lecture
- You can just copy-&-paste others' answers or use AI Tools, but please consider your original stance and opinion too.



Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - > Viterbi decoding

Motivation

- So far: Bag-of-Words (BoW) models
 - Bag = whole document (e.g., Naive Bayes, Vector Space Model)
 - Bag = context of a word (e.g, PPMI and Word2Vec embeddings)
- Natural language: word order matters! (can vary greatly between languages, though)

Bob kills mosquitoes using the book of Hamlet vs. Hamlet kills Bob using the book of mosquitoes

The food tastes good and does not look bad vs. The food tastes bad and does not look good Same words, very different meanings!

Motivation — Example: English

- Fundamental rules word order
 - Subject—Verb—Object (SVO)
 - Adjectives only (immediately) before nouns
 - ...and many more
- "Informal" rules e.g.: order of adjectives
 - $\blacksquare Rule: opinion \rightarrow size \rightarrow physical quality \rightarrow shape \rightarrow age \rightarrow color \rightarrow origin \rightarrow material \rightarrow type \rightarrow purpose$

I saw a beautiful, old, blue, German car.

vs.

I saw a German, blue, beautiful, old car.

Types of Sequence Tasks

• Sequence classification

(Many-to-One, $N \rightarrow 1$)

- Sentiment analysis
- Document categorization



- Sequence labeling/tagging (Many-to-Many, N→N)
 - Part-of-Speech Tagging
 - Named Entity Recognition



Types of Sequence Tasks

• Sequence translation

(Many-to-One, $N \rightarrow M$)

- Machine translation
- Sentence simplification
- Text summarization



• Sequence generation

(One-to-Many, $1 \rightarrow N$)

Image captioning



Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - > Viterbi decoding

Part-of-Speech Tagging

- Part of Speech (also: word class or syntactic category)
 - Each word belongs one or more of these classes
 - English: 8 main parts of speech / word classes (many additional classes and subclasses considered in practice)
- Part-of-Speech (POS) tagging
 - Assign each word in a text a part of speech (duh!)





Penn Treebank Tagset

Base set: 36 POS tags

CC	Coordinating conjunction	and, or
CD	Cardinal number	1, 2, 3, one, two, three
DT	Determiner	the, a, an, any, some
EX	Existential there	
FW	Foreign word	
IN	Preposition / subord. conjunction	in, into, whether, if
JJ	Adjective	cleaner, nice
JJR	Adjective (comparative)	cleaner, nicer
JJS	Adjective (superlative)	cleanes, nicest
LS	List item marker	
MD	Modal	can, could, may
NN	Noun (singular or mass)	machine, computer, air
NNS	Noun (plural)	machines, computers
NNP	Proper noun (singular)	Clementi Mall
NNPS	Proper noun (plural)	Americas
PDT	Predeterminer	all, both, half
POS	Possessive ending	's
PRP	Personal pronoun	him. himself, we

PP\$	Possessive pronoun	her, our, ours	
RB	Adverb	quickly, swiftly	
RBR	Adverb (comparative)	further, greater, more	
RBS	Adverb (superlative)	furthest, greatest, most	
RP	Particle	across, up	
SYM	Symbol	=, +, &	
ТО	to	to	
UH	Interjection	shucks, heck, oops	
VB	Verb (base form)	be, assign, run	
VBD	Verb (past tense)	was, assigned, ran	
VBG	Verb (gerund / present participle)	being, assigning	
VBN	Verb (past participle)	been, assigned	
VBP	Verb (non-3rd pers. sing. present)	am, are	
VBZ	Verb (3rd pers. sing. present)	is	
WDT	wh-determiner	that, which, what	
WP	wh-pronoun	that, which, whom	
WP\$	Possessive wh-pronoun	whose	
WRB	wh-adverb	how, however, why	

Penn Treebank Tagset

Extended set: 12 tags for punctuations and special symbols

#	Pound sign	#
\$	Dollar sign	\$
•	Sentence-final punctuation	.?!
:	Sentence-middle punctuation	:;
,	Comma	,
(Left bracket character	([{ <
)	Right bracket character)]}>
"	Straight double quote	
"	Left open single quote	
"	Left open double quote	
,	Right close single quote	
"	Right close double quote	

Part of Speech — **Two Broad Categories**

• Closed class words

- Small, fixed membership reasonably easy to enumerate
- Generally, short function words that "structure" sentences
- Examples: prepositions, pronouns, participles, determiners, conjunctions, etc.

• Open class words

- Impossible to completely enumerate
- New words continuously being invented, borrowed, etc.
- For most languages: nouns, verbs, adjectives, adverbs

Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - > Viterbi decoding

POS Tagging — Why is it Important?

- Very useful or crucial for many NLP downstream tasks
 - Named Entity recognition (typically comprised of nouns and proper nouns)
 - Information extractions (e.g., verbs indicate relations between entities)
 - Parsing (information of word classes useful before creating parse trees)
 - Speech synthesis/recognition (e.g., noun "DIScount" vs. verb "disCOUNT")
 - Authorship Attribution (e.g., relative frequencies of nouns, verbs, adjectives, etc.)
 - Machine Translation (e.g., reordering of adjectives and nouns)

→ POS tagging: important low-level NLP task

In-Lecture Activity (2 mins)

Quick Quiz



POS Tagging — Why is it Difficult? And How Difficult?

- Our common problem: Ambiguity
 - Many common words have multiple meanings → multiple POS



Often ambiguous, even with additional context

(even humans often don't agree on the correct labeling!)



POS Tagging — Why is it Difficult? And How Difficult?

- POS tagging in English
 - ~85% of word types are unambiguous (e.g., "quickly" is always an adverb, "Alice" is always a proper noun)
 - ~15% of word types are ambiguous but those are quite common!

- → 55-65% of word tokens are ambiguous
 - Ambiguous = 2 or more possible POS tags
 - Results depend on text corpus

In-Lecture Activity (2 mins)

Quick Quiz



In-Lecture Activity (7 mins)

🏃 🏃 🏃 Flying Planes Can Be Dangerous



POS Tagging — Baseline Algorithm

- Most straightforward approach
 - Label each word with its most frequent POS tag
 - Label unknown words as nouns (most common open world class)
- → Result: ~92% accuracy (vs. ~97-98% accuracy for SOTA methods)
 - Doesn't sound so bad, right?
 - 2 main problems:
 - (1) Imbalanced errors
 - High accuracy due to common/frequent unambiguous words (e.g., "the", "a/an", "and", "or")
 - Many of these words also often not that interesting for downstream NLP tasks

(2) Downstream error propagation

■ POS tagging as low-level NLP task → errors quickly propagate up

POS Tagging — Unsupervised Algorithms

- Basic intuition
 - Utilize words with unambiguous POS tags → anchor words
 - Observe patterns to group words into clusters of the same word class
 - Use anchor words to assign clusters (and each containing word) to a POS tag
- Practical considerations
 - No need for hand-labeled text corpora (only lexicon of anchor words required)
 - Poorer performance compared to supervised methods

POS Tagging — Supervised Methods

- Require hand-labeled text corpus
 - Used as input training data for supervised models
 - Challenging for low-resource languages (i.e., languages lacking in large, annotated datasets)
- Popular models (all yielding quite similar results)
 - Hidden Markov Models (HMM)
 - Conditional Random Fields (CRF)
 - Neural sequence models (RNNs, Transformers)
 - Large language models (e.g., BERT)

- Accuracies have reached the human ceiling (i.e., POS taggers as good as human annotators)
- POS tagging considered a solved task for high-resource languages (e.g. English)
- Limitations: low-resource languages and special application domains

In-Lecture Activity (2 mins)

Quick Quiz



Break

Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

• Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - > Viterbi decoding

Markov Chains

• Markov Chain

- Models transitions between a set of <u>states</u> using transition probabilities (captured by a <u>transition matrix</u> A)
- Transition only depends on current state (Markov assumption)
- Sequence: series of transitions
- Example: "Daily Weather"
 - 3 states: *sunny*, *cloudy*, *rainy*

Example question: *"What is the probability of getting a 5 sunny days in a row?"*





Markov Chain → Hidden Markov Models

- Hidden Markov Models (HMM)
 - States are hidden (i.e., not directly observable)
 - Observable variables that depend on the states
- Example: Exercise Routine
 - 3 hidden(!) states: *sunny*, *cloudy*, rainy
 - 2 observed activities: biking, lifting (with the activity depending on the weather)

Example question: "Given that Chris spent the first 3 days lifting and then 3 days biking, what was the most likely weather over the last 6 days?"



HMM — Components

Finite Set of states $S = \{s_1, s_2, ..., s_N\}$

Sequence of states $Q = q_1, q_2, q_3, \dots, q_T$, with $q_t \in S$

Finite set of symbols $V = \{v_1, v_2, ..., v_M\}$

Sequence of observations $O = o_1, o_2, o_3, \dots, o_T$, with $o_t \in V$

Transition probability matrix A $A = \{a_{ij}\}, a_{ij} = P(q_{t+1} = s_j | q_t = s_i)$

Observation / emission probability matrix B $B = \{b_i(o_k)\}, \ b_i(o_k) = P(o_t = v_k | q_t = s_i)$

Initial state distribution π $\pi = \{\pi_i\}, \ \pi_i = P(q_1 \mid s_i)$



HMM — Probabilities (annotated)

Transition probability matrix A

$$A = \{a_{ij}\}, \ a_{ij} = P(q_{t+1} = s_j | \ q_t = s_i) \longleftarrow$$

Probability of transitioning from state s_i to s_j at any time t

$$\sum_{j}^{N} a_{ij} = 1 \ \forall i$$

Observation / emission probability matrix B $B = \{b_i(o_k)\}, \ b_i(o_k) = P(o_t = v_k | q_t = s_i) \leftarrow$

 $\begin{array}{ll} \text{Probability of state } s_i \text{ generating} & & \sum_k^M b_i(o_k) = 1, \ \, \forall k \\ \text{output } v_k \text{ at any time } t & & \sum_k^M b_i(o_k) = 1, \ \, \forall k \end{array}$



HMM — Unrolled Representation

Trellis diagram — graph representation of all possible states and transitions over time



In-Lecture Activity (2 mins)

Quick Quiz



Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

• Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - > Viterbi decoding

POS Tagging with HMMs

- Basic setup
 - Hidden states → POS tags
 - Observations → words
- Example
 - 3 states: {PRP, VBP, NN}



HMM — Core Tasks

(1) Model Learning

Given corresponding state and observation sequences Q and O→ Learn all model parameters, i.e., probabilities A, B and π

(2) Likelihood

Given an HMM $\theta = (A, B, \pi)$

+ an state sequence Q

+ an observation sequence O

→ Calculate the probability $P(Q, O|\theta)$

(3) **Decoding**

Given an HMM $\theta = (A, B, \pi)$ + an observation sequence *O*

→ Find the most likely sequence of states

Training using an annotated dataset

Given 2 POS tag sequences for a sentence, compare which is more likely

Given a sentence, find the most likely POS tags

Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

• Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - > Model Learning
 - Likelihood computation
 - > Viterbi decoding

HMM — Model Learning

• Calculating probabilities using Maximum Likelihood Estimates

$$\pi_i = P(q_1 = s_i) = \frac{Count(\langle S \rangle s_i)}{Count(\langle S \rangle)} \underbrace{ \qquad \text{# sentences starting with state} s_i}_{\text{# sentences}}$$

$$a_{ij} = P(q_{t+1} = s_j | q_t = s_i) = \frac{Count(s_i s_j)}{Count(s_i)} \underbrace{ \overset{\text{--\# occurrences of state}s_i}{\underset{\text{--\# occurrences of state}s_i}} \underbrace{ \overset{\text{--\# occurrences of stat$$

$$b_i(o_k) = P(o_t = v_k | q_t = s_i) = \frac{Count(v_k, s_i)}{Count(s_i)} - \frac{\# \text{ occurrences of observation} v_k \text{ in state} s_i}{F = v_k}$$

HMM — Model Learning — Side Note

• POS tagging using HMM

- Full supervised task → corpus of words labeled with all the correct POS tags
- Fully "visible" Markov Model (we have the state and observation sequences)

"Direct" parameter learning using MLE (we just need simple counts)

- Often in other applications
 - State sequences *Q* are not known
 - Impossible to compute simple counts

Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

• Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - > Viterbi decoding

Likelihood

- Given: HMM $\theta = (A, B, \pi)$ and
 - $O = o_1, o_2, o_3, \dots, o_T$ $Q = q_1, q_2, q_3, \dots, q_T$

• Calculate joint probability $P(O, Q|\theta)$

$$P(O, Q|\theta) = P(O|Q) \cdot P(Q) = \prod_{i=1}^{T} P(o_i|q_i) \cdot P(q_i|q_{i-1}) \qquad \text{with } P(q_1|q_0) = P(q_1)$$

$$emission \\ probabilities \qquad transition \\ probabilities \qquad probabilities \qquad initial state \\ probability \qquad probability \qquad initial state \\ probability \qquad probability \qquad initial state \\ pro$$

Likelihood — Example

$$P(O, Q|\theta) = P(O|Q) \cdot P(Q) = \prod_{i=1}^{T} P(o_i|q_i) \cdot P(q_i|q_{i-1})$$

• Given: sentence O, POS tags Q

 $\left. \begin{array}{l} O = I, like, NLP \\ Q = PRP, VBN, NN \end{array} \right\} \quad P(``I, like, NLP") | PRP-VBN-NN) = ?$

 $P(``I, like, NLP" | PRP - VBN - NN) = P(I|PRP) \cdot P(PRP|\langle S \rangle) \cdot P(like|VBN) \cdot P(VBN|PRP) \cdot P(NLP|NN) \cdot P(NN|VBN)$

All values can be directly taken from *A*, *B*, and π

Likelihood — Example

Visualization using a Trellis diagram \rightarrow just follow the path



$$\begin{split} P(``I, like, NLP") &| PRP - VBN - NN) = \\ P(I|PRP) \cdot P(PRP|\langle S \rangle) \cdot \\ P(like|VBN) \cdot P(VBN|PRP) \cdot \\ P(NLP|NN) \cdot P(NN|VBN) \end{split}$$

Likelihood — Example

Visualization using a Trellis diagram \rightarrow just follow the path



$$\begin{split} P("I, like, NLP" | PRP - VBN - NN) &= \\ P(I|PRP) \cdot P(PRP|\langle S \rangle) \cdot \\ P(like|VBN) \cdot P(VBN|PRP) \cdot \\ P(NLP|NN) \cdot P(NN|VBN) \end{split}$$

 $P(``I, like, NLP" | PRP - VBN - NN) = 0.7 \cdot 0.3 \cdot 0.3 \cdot 0.5 \cdot 0.1 \cdot 0.2$

= 0.00063

In-Lecture Activity (5 mins)

🏃 🏃 🏃 How Complex Could It Be?



Outline

• Overview: Sequence Tasks

• POS Tagging

- What are Parts of Speech?
- Why is this task important and challenging?

• Hidden Markov Models (HMM)

- Basic setup and components
- Core HMM tasks
 - Model Learning
 - Likelihood computation
 - Viterbi decoding

Decoding

- Decoding task
 - Given an HMM $\theta = (A, B, \pi)$ + an observation sequence O
 - Find the most likely sequence of states *Q*



→ **Dynamic Programming** to avoid checking all possible state sequences

Viterbi Algorithm — Toy Example

- Oversimplified setup
 - 3 POS tags: **DT** (determiner), **NN** (noun), **VB** (verb)
 - Let's assume the following HMM

Note: The rows in B do not up to 1 since B does not capture all words, only those we needs.

$$\pi = \begin{bmatrix} 0.8 & 0.2 & 0 \end{bmatrix} \qquad A = \begin{bmatrix} 0 & 0.8 & 0.2 \\ 0 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0 \end{bmatrix} \text{vb} \qquad b = \begin{bmatrix} 0.2 & 0 & 0 & 0 \\ 0.0 & 0.05 & 0.3 & 0.1 \\ 0 & 0.25 & 0.15 & 0.3 \end{bmatrix} \text{vb}$$

• Task: Find the most likely sequence of state (i.e., POS tags) for:

"the fans love the show"





DT NN VB the fans love show $\begin{array}{c} \mathbf{DT} \ \mathbf{NN} \ \mathbf{VB} \\ \pi = \begin{bmatrix} 0.8 \ 0.2 \ 0 \end{bmatrix} \qquad A = \begin{bmatrix} 0 \ 0.8 \ 0.2 \\ 0 \ 0.5 \ 0.5 \\ 0.5 \ 0.5 \ 0 \end{bmatrix} \begin{array}{c} \mathbf{DT} \\ \mathbf{NN} \\ \mathbf{VB} \end{array} \qquad B = \begin{bmatrix} 0.2 \ 0 \ 0 \ 0 \\ 0.0 \ 0.05 \ 0.3 \ 0.1 \\ 0 \ 0.25 \ 0.15 \ 0.3 \end{bmatrix} \begin{array}{c} \mathbf{DT} \\ \mathbf{NN} \\ \mathbf{NN} \\ \mathbf{VB} \end{array}$ Example [DT] DT DT Second word: "fans" • 2 transitions with non-zero probabilities: P(NN|DT)=0.8, P(VB|DT)=0.2 0.8 * 0.2 * 0.8 * 0.05 = 0.16 = <u>0.0064</u> [DT,NN] • "fans" has non-zero emission probabilities for NN and VB: P("fans"|NN)=0.05, P("fans"|VB)=0.25 NN NN 0.2 * 0 • If we would stop here, [DT, VB] would have highest probability • As long there are multiple paths, we have to consider all * 0.2 * 0.25 • We can ignore all paths starting with [DT, DT] 0*0 = 0.008 VB VB VB VB VB [DT,VB] the the show fans love







Viterbi Algorithm

• 2 important question

- How to get the final state sequence with the highest probability?
- How exactly does the Viterbi algorithm reduces complexity?

Backtracking

- During forward pass: remember input path with max probability
- Backtracking: follow paths with max probabilities back to beginning



In-Lecture Activity (2 mins)

Quick Quiz



Viterbi Algorithm — Complexity Analysis



- Let $v_t(3)$ be maximal for the path coming from $v_{t-1}(4)$
- We can ignore all paths coming from $v_{t-1}(j), j \neq 4$
- This holds true for all steps t and states j



Note: Cases where $a_{ij} \cdot b_j(o_t) = 0$ might be an additional convenience but not the main reason for the polynomial complexity

Viterbi Algorithm — The Basic Algorithm

 $v_t(1)$ $v_{t-1}(1)$ -1(2) $v_t(2)$ $v_t(3)$ $v_{t-1}(3)$ $a_{4,3} \cdot b_3(o_t)$ $v_t(4)$ v_{t-1} (4 -1(N) $v_t(N)$ t - 1t

N states

Initialization

$$\begin{aligned} v_1(t) &= \pi_j \cdot b_j(o_1) \\ bt_1(t) &= 0 \end{aligned} \qquad 1 \leq j \leq N \end{aligned}$$

Recursion

$$\begin{aligned} v_t(j) &= \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t) \\ & 1 \le j \le N, \ 1 < t \le T \\ bt_t(j) &= \argmax_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t) \end{aligned}$$

Example for backtrace: $bt_t(3) = 4$ since we get the highest probability for $v_t(3)$ from the path coming from $v_{t-1}(4)$

Viterbi Algorithm — The Basic Algorithm

 $v_t(1)$ $v_{t-1}(1)$ $v_{t-1}(2)$ $v_t(2)$ $v_{t-1}(3)$ $v_t(3)$ $|a_{4,3} \cdot b_3(o_t)|$ $v_t(4)$ $|v_{t-1}|$ $v_{t-1}(N$ $v_t(N)$ t - 1t

N states

Termination (after computing all $v_t(j)$ and $bt_t(j)$)

Probability of most likely path:
$$P^* = \max_{i=1}^N v_T(i)$$

Start of backtrace:
$$q_T^* = \underset{i=1}{\operatorname{argmax}} v_T(i)$$

Viterbi Algorithm — Practical Consideration

• The "usual" problem: Risk of arithmetic underflow

 $v_{1}(t) = \pi_{j} \cdot b_{j}(o_{1})$ $v_{t}(j) = \max_{i=1}^{N} v_{t-1}(i)a_{ij}b_{j}(o_{t})$

Values for $v_t(j)$ become very small as we multiple many (potentially very) small probability values

→ The "usual" solution: Logarithm

$$v_{1}(t) = \log \pi_{j} + \log b_{j}(o_{1})$$
$$v_{t}(j) = \max_{i=1}^{N} v_{t-1}(i) + \log a_{ij} + \log b_{j}(o_{t})$$

Viterbi Algorithm — Python/NumPy Implementation

```
def viterbi(tokens, A, B, PI):
   N, T = A.shape[0], len(tokens)
   M = np.zeros((N, T))
                                           # Reflecting probabilties of trellis
   BT = np.zeros((N, T), dtype=np.int16) # For the Backtracking pointers
   # Initialization
   for s in range(N):
       M[s,0] = PI[s] * B[s, word2index[tokens[0]]]
   # Recursion (with dynamic programming)
   for t in range(1, T):
        for s in range(N):
            new probs = M[:,t-1] * A[:,s] * B[s, word2index[tokens[t]]]
            max idx = np.argmax(new probs)
            M[s,t] = new probs[max idx]
            BT[s,t] = max idx
   # Termination (start backtracking)
   state = np.argmax(M[:,-1])
    state sequence = []
   for i in reversed(range(T)):
        state sequence.append(state)
        state = BT[:,i][state]
   return [ index2tag[idx] for idx in reversed(state sequence) ]
```

Note: This slide is only to show that it does not take much code to implement the Viterbi algorithm.

$$\begin{cases} v_1(t) = \pi_j \cdot b_j(o_1) \qquad bt_1(t) = 0 \\ v_t(j) = \max_{\substack{i=1 \\ i=1}}^N v_{t-1}(i)a_{ij}b_j(o_t) \\ bt_t(j) = \underset{i=1}{\operatorname{argmax}} v_{t-1}(i)a_{ij}b_j(o_t) \\ \end{cases}$$

Viterbi Algorithm — Python/NumPy Implementation

- Using the HMM trained over 25k movie reviews
 - 50 states (POS tags)
 - 83k+ tokens (words, punctuation marks, etc.)

Important: we've cheated here by annotating the reviews using spaCy, not humans!

```
viterbi(['the', 'fans', 'love', 'the', 'show'], A, B, PI)
['DT', 'NNS', 'VBP', 'DT', 'NN']
viterbi(['the', 'fans', 'like', 'the', 'show'], A, B, PI)
['DT', 'NNS', 'IN', 'DT', 'NN']
viterbi(['funny', 'movies', 'are', 'the', 'best'], A, B, PI)
['JJ', 'NNS', 'VBP', 'DT', 'JJS']
viterbi(['i', 'like', 'watching', 'comedies'], A, B, PI)
['PRP', 'VBP', 'VBG', 'NNS']
```

Summary

• Sequences

- A primary form of natural language data with many applications (a sentence is sequence of words; sequence captures meaning → BoW model intrinsically limited)
- Many sequence tasks in NLP
- Focus of this lecture: sequence labeling
 - POS tagging as very fundamental sequence labeling task
 - Different approaches, incl. Hidden Markov Models (HMM)
- Next lecture: encoder-decoder architecture
 - Neural network-based architecture
 - Applicable to all sequence tasks

Outlook for Next Week: Encoder–Decoder



GIF Credits Aditya Shirsath @ Medium

Pre-Lecture Activity for Next Week

• Assigned Task

Post a 1-2 sentence answer to the following question in the [Pre-Lecture] discussion

"What is an encoder? What is a decoder? What are we trying to encode/decode anyways?"

Side notes:

- This task is meant as a warm-up to provide some context for the next lecture
- No worries if you get lost; we will talk about this in the next lecture
- You can just copy-&-paste others' answers but this won't help you learn better