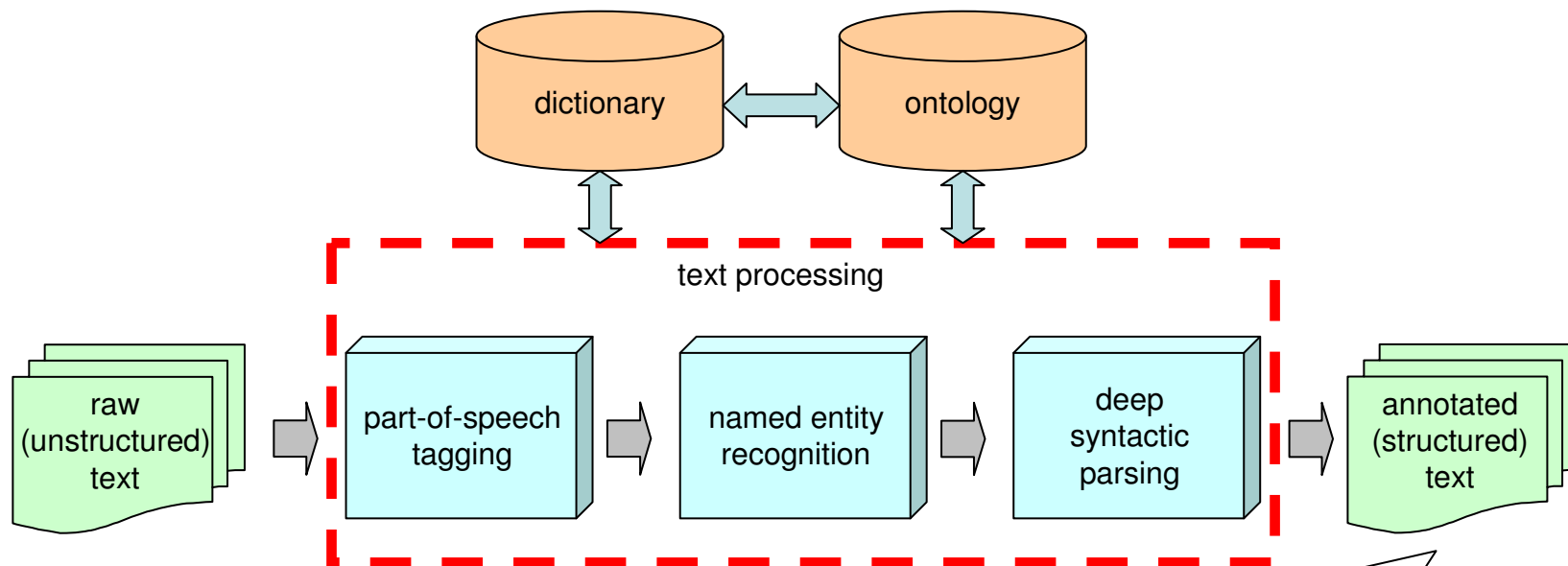


Sentence splitting, tokenization,
language modeling, and part-of-
speech tagging with hidden
Markov models

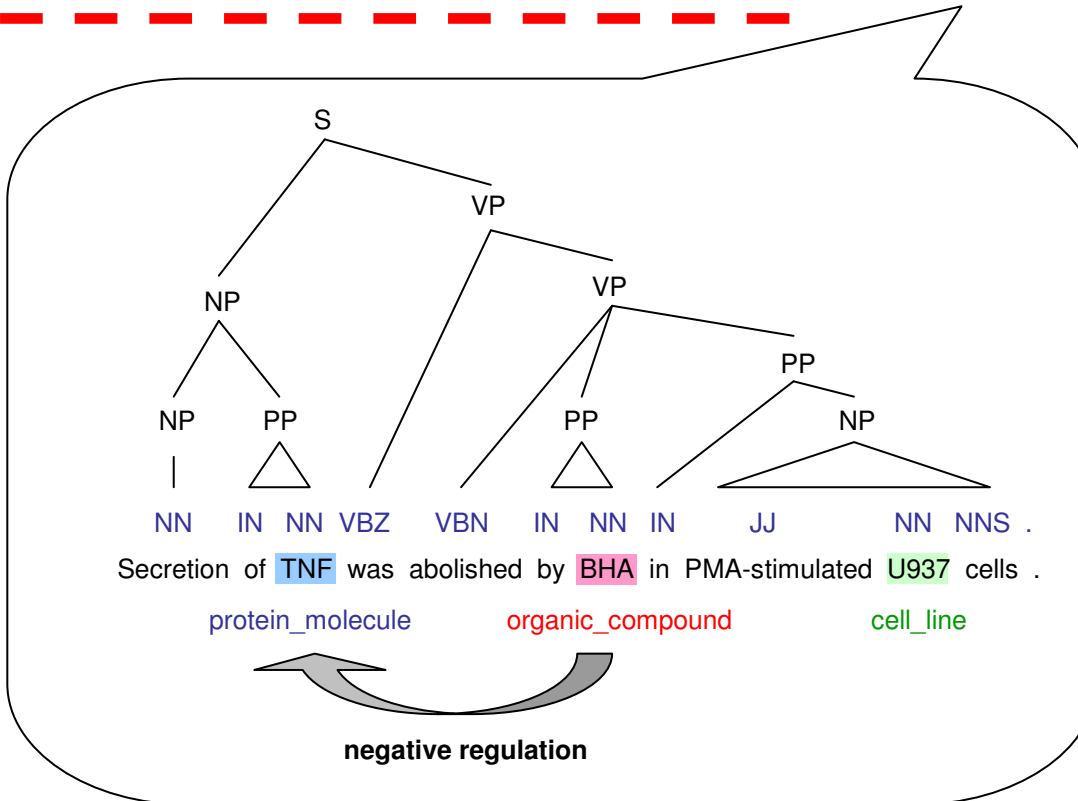
Yoshimasa Tsuruoka

Topics

- Sentence splitting
- Tokenization
- Maximum likelihood estimation (MLE)
- Language models
 - Unigram
 - Bigram
 - Smoothing
- Hidden Markov models (HMMs)
 - Part-of-speech tagging
 - Viterbi algorithm



.....
 ... Secretion of TNF was abolished by BHA in PMA-stimulated U937 cells.



Basic Steps of Natural Language Processing

- Sentence splitting
- Tokenization
- Part-of-speech tagging
- Shallow parsing
- Named entity recognition
- Syntactic parsing
- (Semantic Role Labeling)

Sentence splitting

Current immunosuppression protocols to prevent lung transplant rejection reduce pro-inflammatory and T-helper type 1 (Th1) cytokines. However, Th1 T-cell pro-inflammatory cytokine production is important in host defense against bacterial infection in the lungs. Excessive immunosuppression of Th1 T-cell pro-inflammatory cytokines leaves patients susceptible to infection.



Current immunosuppression protocols to prevent lung transplant rejection reduce pro-inflammatory and T-helper type 1 (Th1) cytokines.

However, Th1 T-cell pro-inflammatory cytokine production is important in host defense against bacterial infection in the lungs.

Excessive immunosuppression of Th1 T-cell pro-inflammatory cytokines leaves patients susceptible to infection.

A heuristic rule for sentence splitting

sentence boundary

= period + space(s) + capital letter

Regular expression in Perl

```
s/\. +([A-Z])\.\n$/g;
```

Errors

IL-33 is known to induce the production of Th2-associated cytokines (e.g. IL-5 and IL-13).



IL-33 is known to induce the production of Th2-associated cytokines (e.g.

IL-5 and IL-13).

- Two solutions:
 - Add more rules to handle exceptions
 - Machine learning

Adding rules to handle exceptions

```
#!/usr/bin/perl

while(<STDIN>) {
    $_ =~ s/([\.\?]) +([\(0-9a-zA-Z)]/$1\n$2/g;

    $_ =~ s/(\WDr\.)\n/$1 /g;
    $_ =~ s/(\WMr\.)\n/$1 /g;
    $_ =~ s/(\WMs\.)\n/$1 /g;
    $_ =~ s/(\WMrs\.)\n/$1 /g;
    $_ =~ s/(\Wvs\.)\n/$1 /g;
    $_ =~ s/(\Wa\.m\.)\n/$1 /g;
    $_ =~ s/(\Wp\.m\.)\n/$1 /g;
    $_ =~ s/(\Wi\.e\.)\n/$1 /g;
    $_ =~ s/(\We\.g\.)\n/$1 /g;

    print;
}
```


Tokenization

“We’re heading into a recession.”



“	We	're	heading	into	a	recession	.	”
----------	-----------	------------	----------------	-------------	----------	------------------	----------	----------

- Tokenizing general English sentences is relatively straightforward.
- Use spaces as the boundaries
- Use some heuristics to handle exceptions

Exceptions

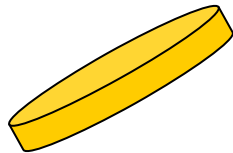
- separate possessive endings or abbreviated forms from preceding words:
 - Mary's → Mary 's
 - Mary's → Mary is
 - Mary's → Mary has
- separate punctuation marks and quotes from words :
 - Mary. → Mary .
 - “new” → “ new ”

Tokenization

- Tokenizer.sed: a simple script in *sed*
 - <http://www.cis.upenn.edu/~treebank/tokenization.html>
- Undesirable tokenization
 - original: “1,25(OH)2D3”
 - tokenized: “1 , 25 (OH) 2D3”
- Tokenization for biomedical text
 - Not straight-forward
 - Needs dictionary? Machine learning?

Maximum likelihood estimation

- Learning parameters from data
- Coin flipping example
 - We want to know the probability that a biased coin comes up heads.
 - Data: {H, T, H, H, T, H, T, T, H, H, T}



$$P(\text{Head}) = \frac{6}{11}$$

Why?

Maximum likelihood estimation

- Likelihood

- Data: {H, T, H, H, T, H, T, T, H, H, T}

$$\begin{aligned} P(\text{Data}) &= \theta \cdot (1 - \theta) \cdot \theta \cdot \theta \cdot (1 - \theta) \cdot \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta \cdot \theta \cdot (1 - \theta) \\ &= \theta^6 (1 - \theta)^5 \end{aligned}$$

- Maximize the (log) likelihood

$$\log P(\text{Data}) = 6 \log \theta + 5 \log(1 - \theta)$$

$$\frac{d \log P(\text{Data})}{d\theta} = \frac{6}{\theta} - \frac{5}{1 - \theta} = 0$$



$$\theta = \frac{6}{11}$$

Language models

- A language model assigns a probability to a word sequence $P(w_1 \dots w_n)$
- Statistical machine translation

English sentence

$$P(e|f) = \frac{P(f|e)P(e)}{P(f)} \propto P(f|e)P(e)$$

French sentence

Language model for English

- Noisy channel models: machine translation, part-of-speech tagging, speech recognition, optical character recognition, etc.

Language modeling

$w_1 \dots w_n$	$P(w_1 \dots w_n)$
He opened the window .	0.0000458
She opened the window .	0.0000723
John was hit .	0.0000244
John hit was .	0.0000002
:	:

- How do you compute $P(w_1 \dots w_n)$?
 - ➡ Learn from data (a corpus)

Estimate probabilities from a corpus

Corpus

Humpty Dumpty sat on a wall .
Humpty Dumpty had a great fall .

$$P(\text{Humpty Dumpty sat on a wall .}) = \frac{1}{2}$$

$$P(\text{Humpty Dumpty had a great fall .}) = \frac{1}{2}$$

$$P(\text{Humpty Dumpty had a fall .}) = 0$$

Sequence to words

- Let's decompose the sequence probability into probabilities of individual words

$$\begin{aligned} P(w_1 \dots w_n) &= P(w_1)P(w_2 \dots w_n \mid w_1) \quad \leftarrow \text{Chain rule } P(x, y) = P(x)P(y \mid x) \\ &= P(w_1)P(w_2 \mid w_1)P(w_3 \dots w_n \mid w_1 w_2) \\ &= P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1 w_2)P(w_4 \dots w_n \mid w_1 w_2 w_3) \\ &= \dots \\ &= \prod_{i=1}^n P(w_i \mid w_1 \dots w_{i-1}) \end{aligned}$$

Unigram model

- Ignore the context completely

$$P(w_1 \dots w_n) = \prod_{i=1}^n P(w_i \mid w_1 \dots w_{i-1})$$
$$\approx \prod_{i=1}^n P(w_i)$$

- Example

$$P(\text{Humpty Dumpty sat on a wall.})$$
$$\approx P(\text{Humpty})P(\text{Dumpty})P(\text{sat})P(\text{on})P(\text{a})P(\text{wall})P(.)$$

MLE for unigram models

Humpty Dumpty sat on a wall .
Humpty Dumpty had a great fall .

- Count the frequencies

$$P(w_k) = \frac{C(w_k)}{\sum_w C(w_k)}$$

$$P(\text{Humpty}) = \frac{2}{14}$$

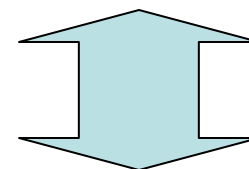
$$P(\text{on}) = \frac{1}{14}$$

Unigram model

- Problem

$$P(\text{Humpty Dumpty sat on a wall.}) \\ \approx P(\text{Humpty})P(\text{Dumpty})P(\text{sat})P(\text{on})P(\text{a})P(\text{wall})P(.)$$

$$P(\text{Dumpty Humpty sat on a wall.}) \\ \approx P(\text{Humpty})P(\text{Dumpty})P(\text{sat})P(\text{on})P(\text{a})P(\text{wall})P(.)$$



- Word order is not considered.

Bigram model

- 1st order Markov assumption

$$P(w_1 \dots w_n) = \prod_{i=1}^n P(w_i | w_1 \dots w_{i-1})$$
$$\approx \prod_{i=1}^n P(w_i | w_{i-1})$$

- Example

$P(\text{Humpty Dumpty sat on a wall.})$

$$\approx P(\text{Humpty} | \langle s \rangle) P(\text{Dumpty} | \text{Humpty}) P(\text{sat} | \text{Dumpty})$$
$$P(\text{on} | \text{sat}) P(\text{a} | \text{on}) P(\text{wall} | \text{a}) P(. | \text{wall})$$

MLE for bigram models

Humpty Dumpty sat on a wall .
Humpty Dumpty had a great fall .

- Count frequencies

$$P(w_k | w_{k-1}) = \frac{C(w_{k-1}w_k)}{C(w_{k-1})}$$

$$P(\text{sat} | \text{Dumpty}) = \frac{1}{2}$$

$$P(\text{had} | \text{Dumpty}) = \frac{1}{2}$$

N-gram models

- $n-1$ order Markov assumption

$$P(w_1 \dots w_n) = \prod_{i=1}^n P(w_i \mid w_1 \dots w_{i-1})$$
$$\approx \prod_{i=1}^n P(w_i \mid w_{i-n+1} \dots w_{i-1})$$

- Modeling with a large n
 - Pros: rich contextual information
 - Cons: sparseness (unreliable probability estimation)

Web 1T 5-gram data

- Released from Google
 - <http://www ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC2006T13>
 - 1 trillion word tokens of text from Web pages
- Examples

serve as the incoming	92
serve as the incubator	99
serve as the independent	794
serve as the index	223
serve as the indication	72
serve as the indicator	120
:	:

Smoothing

- Problems with the maximum likelihood estimation method
 - Events that are not observed in the training corpus are given zero probabilities
 - Unreliable estimates when the counts are small

Add-One Smoothing

Humpty Dumpty sat on a wall .
Humpty Dumpty had a great fall .

- V : the size of the vocabulary

$$P(w_k | w_{k-1}) = \frac{C(w_{k-1}w_k) + 1}{C(w_{k-1}) + V}$$

$$P(\text{sat} | \text{Dumpty}) = \frac{1+1}{2+10}$$

Part-of-speech tagging

Paul Krugman, a professor at Princeton University, was
NNP NNP , DT NN IN NNP NNP , VBD
awarded the Nobel Prize in Economics on Monday.
VBN DT NNP NNP IN NNP IN NNP

- Assigns a part-of-speech tag to each word in the sentence.
- Part-of-speech tags

NN: Noun

NNP: Proper noun

DT: Determiner

:

IN: Preposition

VBD: Verb, past tense

VBN: Verb, past participle

:

Part-of-speech tagging is not easy

- Parts-of-speech are often ambiguous

I have to go there.
verb

I had a go at it.
noun

- We need to look at the context
- But how?

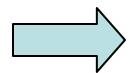
Writing rules for part-of-speech tagging

I have to go there.
verb

I had a go at it.
noun

- If the previous word is “to”, then it’s a verb.
- If the previous word is “a”, then it’s a noun.
- If the next word is ...

:



Writing rules manually is impossible

Learning from examples

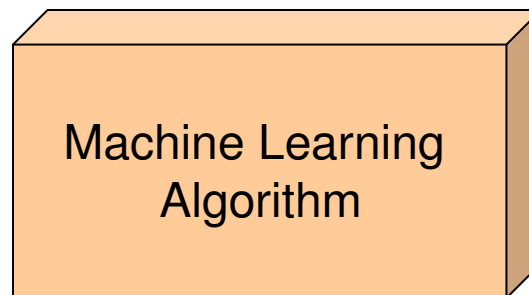
The involvement of ion channels in B and T lymphocyte activation is
DT NN IN NN NNS IN NN CC NN NN NN VBZ
supported by many reports of changes in ion fluxes and membrane
VBN IN JJ NNS IN NNS IN NN NNS CC NN

.....

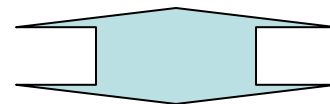
.....

Unseen text

We demonstrate
that ...



We demonstrate
PRP VBP
that ...
IN



training

Part-of-speech tagging

- Input:
 - Word sequence $W = w_1 \dots w_n$
- Output (prediction):
 - The most probable tag sequence $T = t_1 \dots t_n$

$$\begin{aligned} \hat{T} &= \operatorname{argmax}_{T \in \tau} P(T | W) \\ &= \operatorname{argmax}_{T \in \tau} \frac{P(T)P(W | T)}{P(W)} \quad \leftarrow \text{constant} \\ &= \operatorname{argmax}_{T \in \tau} P(T)P(W | T) \end{aligned}$$

Assumptions

$$\begin{aligned}P(W | T) &= P(w_1 \dots w_n | t_1 \dots t_n) \\&= P(w_1 | t_1 \dots t_n) P(w_2 \dots w_n | w_1 t_1 \dots t_n) \\&= P(w_1 | t_1 \dots t_n) P(w_2 | w_1 t_1 \dots t_n) P(w_3 \dots w_n | w_1 w_2 t_1 \dots t_n) \\&= \dots \\&= \prod_{i=1}^n P(w_i | w_1 \dots w_{i-1} t_1 \dots t_n)\end{aligned}$$

- Assume the word emission probability depends only on its tag

$$P(W | T) = \prod_{i=1}^n P(w_i | t_i)$$

Part-of-speech tagging with Hidden Markov Models

- Assume the first-order Markov assumption for tags

$$P(T) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

- Then, the most probably tag sequence can be expressed as

$$\hat{T} = \operatorname{argmax}_{T \in \tau} P(T)P(W | T)$$

$$= \operatorname{argmax}_{T \in \tau} \prod_{i=1}^n P(t_i | t_{i-1})P(w_i | t_i)$$

transition probability

emission probability

Learning from training data

The involvement of ion channels in B and T lymphocyte activation is

DT NN IN NN NNS IN NN CC NN NN NN VBZ

supported by many reports of changes in ion fluxes and membrane

VBN IN JJ NNS IN NNS IN NN NNS CC NN

.....
.....



Maximum likelihood estimation

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

emission probability $P(w_i | t_i) = \frac{C(w_i t_i)}{C(t_i)}$

Examples from the Wall Street Journal corpus (1)

$$P(\text{NN} | \text{DT}) = 0.47$$

$$P(\text{JJ} | \text{DT}) = 0.22$$

$$P(\text{NNP} | \text{DT}) = 0.11$$

$$P(\text{NNS} | \text{DT}) = 0.07$$

$$P(\text{CD} | \text{DT}) = 0.02$$

$$P(\text{IN} | \text{NN}) = 0.25$$

$$P(\text{NN} | \text{NN}) = 0.12$$

$$P(, | \text{NN}) = 0.11$$

$$P(. | \text{NN}) = 0.08$$

$$P(\text{VBD} | \text{NN}) = 0.05$$

POS Tag	Description	Examples
DT	determiner	a, the
NN	noun, singular or mass	company
NNS	noun plural	companies
NNP	proper nouns, singular	Nasdaq
IN	preposition/subordinating conjunction	in, of, like
VBD	verb, past tense	took
:	:	:

Examples from the Wall Street Journal corpus (2)

$$P(\% \mid \text{NN}) = 0.037$$

$$P(\text{company} \mid \text{NN}) = 0.019$$

$$P(\text{year} \mid \text{NN}) = 0.0167$$

$$P(\text{market} \mid \text{NN}) = 0.0155$$

$$P(\text{share} \mid \text{NN}) = 0.0107$$

$$P(\text{said} \mid \text{VBD}) = 0.188$$

$$P(\text{was} \mid \text{VBD}) = 0.131$$

$$P(\text{were} \mid \text{VBD}) = 0.064$$

$$P(\text{had} \mid \text{VBD}) = 0.056$$

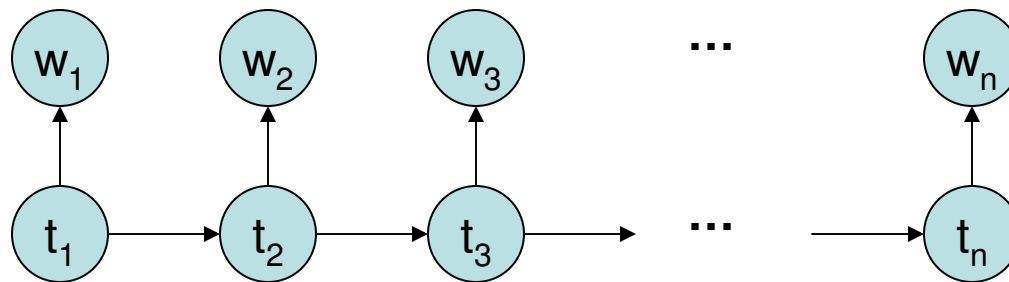
$$P(\text{rose} \mid \text{VBD}) = 0.022$$

POS Tag	Description	Examples
DT	determiner	a, the
NN	noun, singular or mass	company
NNS	noun plural	companies
NNP	proper nouns, singular	Nasdaq
IN	preposition/subordinating conjunction	in, of, like
VBD	verb, past tense	took
:	:	:

Part-of-speech tagging

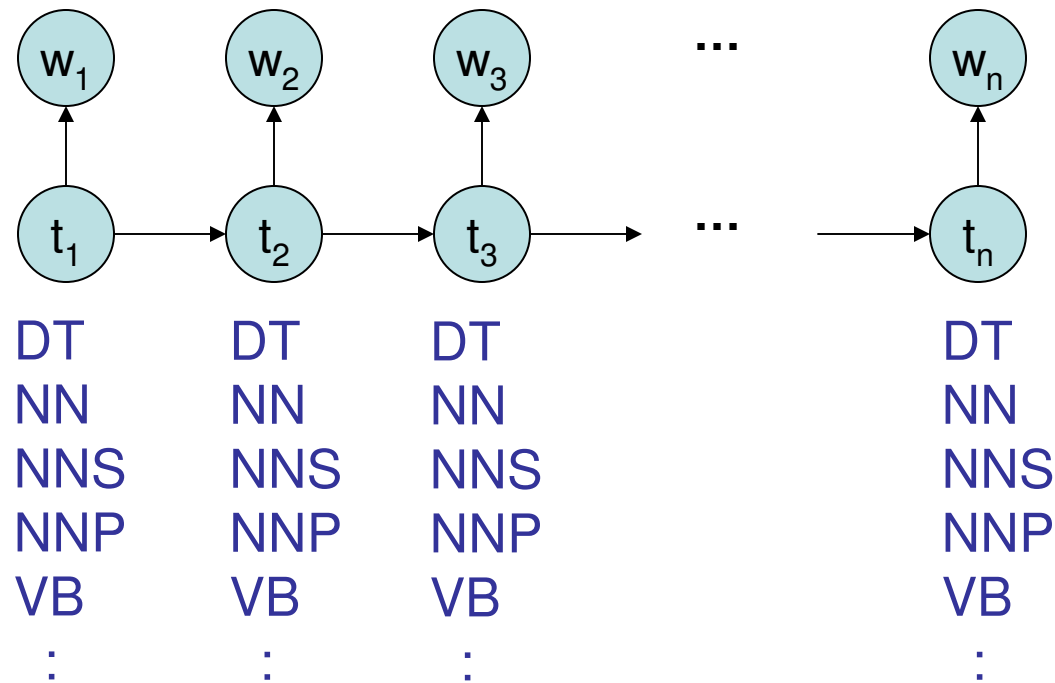
- Input:
 - Word sequence $W = w_1 \dots w_n$
- Output (prediction):
 - The most probable tag sequence $T = t_1 \dots t_n$

$$\hat{T} = \operatorname{argmax}_{T \in \tau} \prod_{i=1}^n P(t_i | t_{i-1}) P(w_i | t_i)$$



Finding the best sequence

- Naive approach:
 - Enumerate all possible tag sequences with their probabilities and select the one that gives the highest probability



→ exponential in the length of the sentence

Finding the best sequence

- if you write them down...

$$\begin{aligned} &P(\text{DT} | \langle s \rangle)P(w_1 | \text{DT})P(\text{DT} | \text{DT})P(w_2 | \text{DT}) \wedge P(\text{DT} | \text{DT})P(w_n | \text{DT}) \\ &P(\text{NN} | \langle s \rangle)P(w_1 | \text{NN})P(\text{DT} | \text{DT})P(w_2 | \text{DT}) \wedge P(\text{DT} | \text{DT})P(w_n | \text{DT}) \\ &P(\text{NNS} | \langle s \rangle)P(w_1 | \text{NNS})P(\text{DT} | \text{DT})P(w_2 | \text{DT}) \wedge P(\text{DT} | \text{DT})P(w_n | \text{DT}) \\ &\quad \vdots \\ &P(\text{DT} | \langle s \rangle)P(w_1 | \text{DT})P(\text{NN} | \text{DT})P(w_2 | \text{NN}) \wedge P(\text{DT} | \text{DT})P(w_n | \text{DT}) \\ &P(\text{NN} | \langle s \rangle)P(w_1 | \text{NN})P(\text{NN} | \text{NN})P(w_2 | \text{NN}) \wedge P(\text{DT} | \text{DT})P(w_n | \text{DT}) \\ &P(\text{NNS} | \langle s \rangle)P(w_1 | \text{NNS})P(\text{NN} | \text{NNS})P(w_2 | \text{NN}) \wedge P(\text{DT} | \text{DT})P(w_n | \text{DT}) \end{aligned}$$

A lot of redundancy in computation



There should be a way to do it more efficiently

Viterbi algorithm (1)

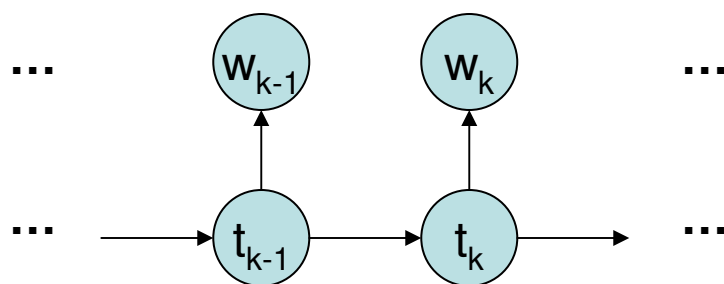
- Let $V_k(t_k)$ be the probability of the best sequence for $w_1 \dots w_k$ ending with the tag t_k

$$\begin{aligned} V_k(t_k) &\equiv \max_{t_1, t_2, \dots, t_{k-1}} \prod_{i=1}^k P(t_i | t_{i-1}) P(w_i | t_i) \\ &= \max_{t_{k-1}} \max_{t_1, t_2, \dots, t_{k-2}} \prod_{i=1}^k P(t_i | t_{i-1}) P(w_i | t_i) \\ &= \max_{t_{k-1}} P(t_k | t_{k-1}) P(w_k | t_k) \max_{t_1, t_2, \dots, t_{k-2}} \prod_{i=1}^{k-1} P(t_i | t_{i-1}) P(w_i | t_i) \\ &= \max_{t_{k-1}} P(t_k | t_{k-1}) P(w_k | t_k) \underline{V_{k-1}(t_{k-1})} \end{aligned}$$

We can compute them in a recursive manner!

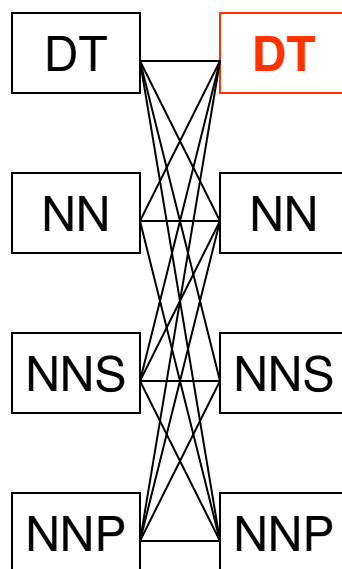
Viterbi algorithm (2)

$$V_k(t_k) = \max_{t_{k-1}} P(t_k | t_{k-1}) P(w_k | t_k) V_{k-1}(t_{k-1})$$



Beginning of the sentence

$$V_0(< s >) = 1$$

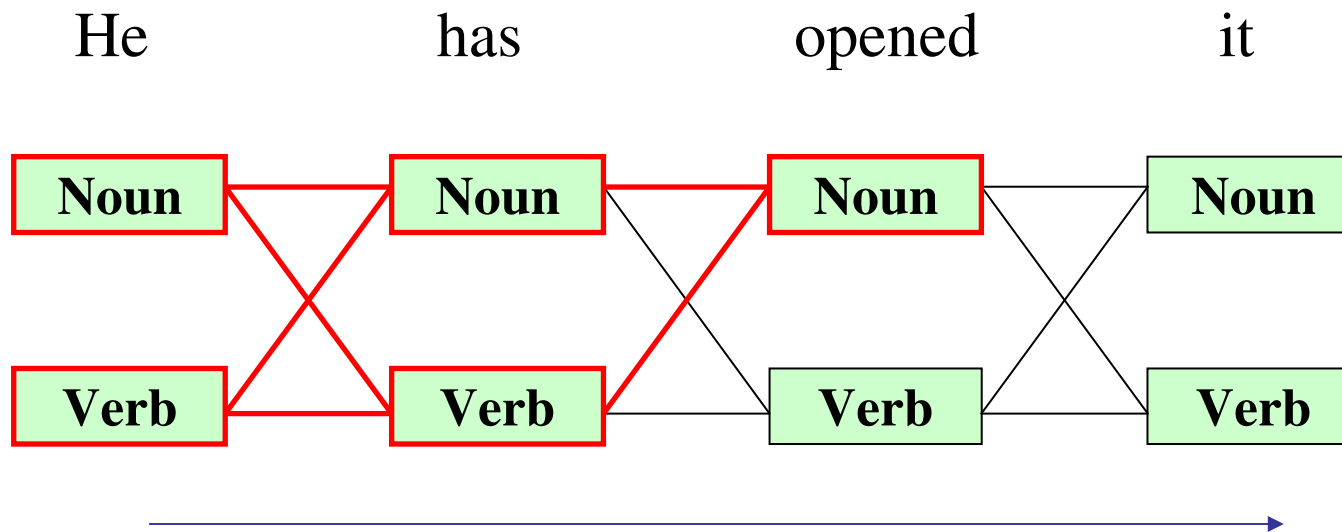


$$\begin{aligned}
 &P(\text{DT} | \text{DT}) P(w_k | \text{DT}) V_{k-1}(\text{DT}) \\
 &P(\text{DT} | \text{NN}) P(w_k | \text{DT}) V_{k-1}(\text{NN}) \\
 &P(\text{DT} | \text{NNS}) P(w_k | \text{DT}) V_{k-1}(\text{NNS}) \\
 &P(\text{DT} | \text{NNP}) P(w_k | \text{DT}) V_{k-1}(\text{NNP}) \\
 &\quad \vdots
 \end{aligned}$$

Viterbi algorithm (3)

- Left to right
 - Save backward pointers to recover the best tag sequence

Complexity: $(\text{number of tags})^2 \times (\text{length})$



Topics

- Sentence splitting
- Tokenization
- Maximum likelihood estimation (MLE)
- Language models
 - Unigram
 - Bigram
 - Smoothing
- Hidden Markov models (HMMs)
 - Part-of-speech tagging
 - Viterbi algorithm