

Natural Language Understanding and Generation

Introduction to Natural Language Processing

Challenges

Syntax & Semantics

- L = Words + rules + exceptions..
- Ambiguity at all levels..

Variety

- We speak different languages..
- And language is a cultural entity..
- So they are not equivalent..
- Highly systematic but also complex..
- Keeps changing.. New words, New rules and New exceptions..

Source :

Electronic texts / Printed texts / Acoustic Speech Signal..

- they are noisy..
- Language looks obvious to us.. But it is a Big Deal

Lexical analysis is dividing the whole chunk of text into paragraphs, sentences, and words. Concerns with listing of words and categorizing them word and word meanings - the lexicon

Lexicon of a language means the collection of words and phrases in a language. In linguistics, a lexicon is a language's inventory of lexemes

friendly or beautyship → Lexically Incorrect

friendship and beautiful is correct → Lexically Correct

Semantic knowledge is the aspect of language knowledge that involves word meanings/vocabulary.

Concerns what words mean and how these meanings – combine in sentences to form sentence meanings.

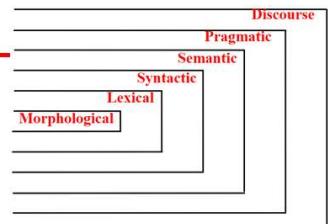
Semantics concerns the meaning of linguistic expressions independently of the context in which they occur.

"Green ideas have large noses" → Syntactically Correct But Semantically Incorrect

Different Levels of Language Analysis

Morphological Analysis

- Lexical Analysis
- Syntactic Analysis
- Semantics Analysis
- Pragmatic Analysis
- Discourse Knowledge



Types of Ambiguities:

Structural Ambiguities

- Namrata thinks she understands me.
- She thinks Namrata understands me.
- Visiting relatives can be nuisance. (two meanings)

Grammatical Ambiguities

- I (feminine or masculine) go.
- Can- Noun = container, Can – Modal(auxiliary verb), Can-verb = to can means to pack

Lexical Ambiguities:

- Polysemy Ex: "understand" (I get it)
- Ex: going to the bank/savings account in a bank
(BUILDIN <-> ORGANIZATION)
- Homonymy Ex: Bank= river, financial bank

Morphology concerns how words are constructed from basic meaning units called morphemes. A morpheme is the primitive unit of meaning in a language

FireHouse = Fire + House each morpheme is Free morphemes

has unique meaning morphemes

cats → morphological parser → cat + s → No meaning bound
↳ some meaning morphemes

Syntax analysis checks the text for meaningfulness comparing to the rules of formal language. Concerns how words can be put together to form correct sentences and determines what structural role each word plays in the sentence and what phrases are subparts of other phrases

"Large have green ideas nose" → Lexically Correct But Syntactically Incorrect

Pragmatics is meaning in use / context. Concerns how sentences are used in different situations how use affects the interpretation of sentence and whether it produces a useful meaning.

"She cuts banana with a pen" → Semantically Correct But Pragmatically Incorrect

Discourse analysis concerns how the immediately preceding sentences affect the interpretation of next sentence

"NLP is a very interesting subject. IT is liked by students. They generally do well in exams." → Pragmatically Correct But Incorrect at a Discourse Level

- Morphological** analysis: This is the analysis of the structure and formation of words. For example, the word 'unhappy' consists of two morphemes: the prefix 'un-' and the root 'happy'. The prefix 'un-' changes the meaning of the root 'happy' to its opposite. A small example of morphological analysis is to identify the morphemes and their functions in a word like 'reopened'. A counterexample to the statement "Every word consists of one or more morphemes" is the word 'a', which is not a morpheme by itself, but a phonetic reduction of the morpheme 'an'.
- Lexical** analysis: This is the analysis of the meaning and usage of words and phrases. For example, the word 'bank' can have different meanings depending on the context: it can mean a financial institution, a slope of land, or a set of similar things. A small example of lexical analysis is to explain the difference between the words 'affect' and 'effect'. A counterexample to the statement "Every word has a fixed meaning" is the word 'bat', which can mean a flying mammal, a wooden club, or a verb meaning to strike something.
- Syntactic** analysis: This is the analysis of the structure and rules of sentences and clauses. For example, a sentence consists of a subject and a predicate, and a clause consists of a subject, a verb, and an object or complement. A small example of syntactic analysis is to identify the parts of speech and the functions of the words in a sentence like 'She gave him a book'. A counterexample to the statement "Every sentence has a subject and a predicate" is the sentence 'Go!', which has only an implied subject and a predicate.
- Semantic** analysis: This is the analysis of the meaning and interpretation of sentences and texts. For example, a sentence can have different meanings depending on the context, the tone, or the intention of the speaker or writer. A small example of semantic analysis is to explain the difference between the sentences 'He is cold' and 'He is a cold person'. A counterexample to the statement "Every sentence has a literal meaning" is the sentence 'It's raining cats and dogs', which has a figurative meaning that implies heavy rain, not animals falling from the sky.
- Pragmatic** analysis: This is the analysis of the use and effect of language in specific situations and purposes. For example, language can be used to persuade, inform, entertain, or request. A small example of pragmatic analysis is to explain the difference between the utterances 'Can you pass me the salt?' and 'Pass me the salt'. A counterexample to the statement "Every utterance has a direct meaning" is the utterance 'Can you open the window?', which has an indirect meaning that requests the hearer to open the window, not a question about their ability.
- Discourse** analysis: This is the analysis of the structure and coherence of larger units of language, such as paragraphs, conversations, or texts. For example, discourse analysis can examine how different speakers or writers organize their ideas, how they relate to each other, and how they achieve their goals. A small example of discourse analysis is to identify the topic, the main points, and the transitions in a paragraph. A counterexample to the statement "Every paragraph has a topic sentence" is the paragraph 'I like apples. They are sweet and crunchy. You can make pies and jams with them. They are good for your health.', which has no clear topic sentence, but several supporting sentences.

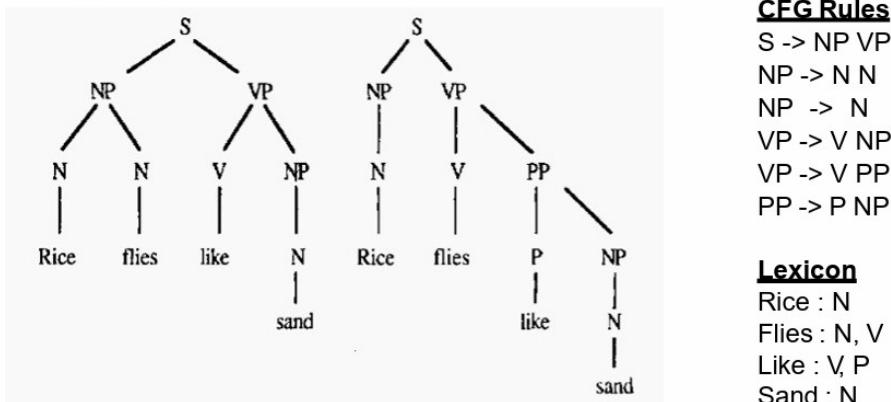


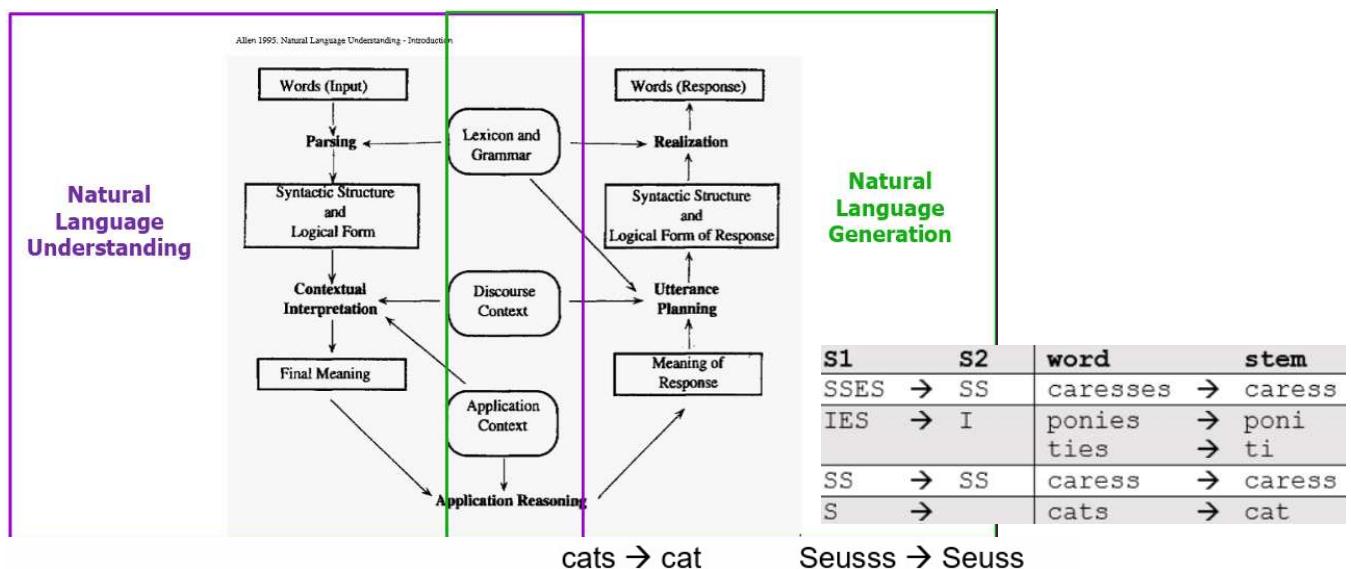
figure 1.4 Two structural representations of *Rice flies like sand*.

NP – Noun Phrases

VP – Verb Phrases

PP – Prepositional Phrases

1.4 Two structural representations of "Rice flies like sand".



Stemming is a crude heuristic process that chops off the ends of words, often includes the removal of derivational affixes

"Seusss", "cat", "hat", "different", "other", "cats", "Seuss", "perfectionist", "always", "organized", "Being", "firm", "believer", "democracy", "seuss", "ardent", "fan", "democratic" party", "nation"

break, breaks, broke, broken and breaking are forms of the same **lexeme**, with break as the **lemma** by which they are **indexed**.

Lemma is the canonical form, dictionary form, or citation form of a set of word forms.

Lemma: same stem, part of speech, rough word sense

Lemmatization is the task of determining that two words have the same root despite their surface differences.

cat and **cats** = same lemma cat
am, are, and is : lemma be

The word "**walk**" is the base form for the word "**walking**", and hence this is matched in both stemming and lemmatization.

"meeting" : base form of a noun
: or a form of a verb ("**to meet**")
e.g., "in our last meeting" or "We are meeting again tomorrow"

Wordform: the full inflected surface form
cat and **cats** = different wordforms

"Seuss", "cat", "hat", "different", "other", "cat", "Seuss", "perfectionist", "always", "organized", "Being", "firm", "believer", "democracy", "seuss", "ardent", "fan", "democratic" party", "nation"

N-Gram Language Modelling

Simplifying assumption:

limit history of fixed number of words, n

$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{that})$
or

$P(\text{the} \mid \text{its water is so transparent that}) \approx P(\text{the} \mid \text{transparent that})$

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i \mid w_{i-k} \dots w_{i-1})$$

$$P(w_n \mid w_{n-1}) = \frac{C(w_{n-1} w_n)}{C(w_{n-1})}$$

<s> I am Sam </s>
 <s> Sam I am </s>
 <s> I do not like green eggs and ham </s>

In other words, we approximate each component in the product

$$P(\text{I} \mid \text{<s>}) = \frac{2}{3} = .67 \quad P(\text{Sam} \mid \text{<s>}) = \frac{1}{3} = .33 \quad P(\text{am} \mid \text{I}) = \frac{2}{3} = .67$$

$$P(w_i \mid w_1 w_2 \dots w_{i-1}) \approx P(w_i \mid w_{i-k} \dots w_{i-1}) \quad P(\text{</s>} \mid \text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam} \mid \text{am}) = \frac{1}{2} = .5 \quad P(\text{do} \mid \text{I}) = \frac{1}{3} = .33$$

N=1 → Unigram language model

$$P(w_1 w_2 \dots w_n) \approx \prod_i^n P(w_i) \quad P(w_{1:n}) \approx \prod_{k=1}^{n-1} P(w_k \mid w_{k-1})$$

$$P(\text{"its water is so transparent"}) = P(\text{its}) \times P(\text{water}) \times P(\text{is}) \times P(\text{so}) \times P(\text{transparent})$$

N=2 → Bigram language model

$$P(\text{"I want to eat healthy food"}) =$$

$$P(\text{I} \mid \text{<start>}) * P(\text{want} \mid \text{I}) * P(\text{to} \mid \text{want}) * P(\text{eat} \mid \text{to}) * P(\text{healthy} \mid \text{eat}) * P(\text{food} \mid \text{healthy}) * P(\text{<end>} \mid \text{food})$$

Perplexity is a measure of how well a language model predicts a sequence of words. It is defined as the exponentiated average negative log-likelihood of a sequence, calculated with exponent base $e^{\frac{1}{N}}$. In other words, perplexity is the inverse probability of the sequence, normalized by its length. A lower perplexity means a higher probability, and a higher perplexity means a lower probability.

Perplexity tells us how surprised or uncertain a language model is about the next word in the sequence. A good language model should assign high probabilities to well-formed and coherent sentences, and low probabilities to ill-formed and incoherent sentences. Therefore, a good language model should have a low perplexity on the test data. Perplexity is one way to evaluate the quality and performance of language models in NLP ².

$$\text{perplexity}(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i \mid w_1 \dots w_{i-1})}} \quad \text{unigram} \quad \text{perplexity}(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i)}} \quad \text{bigram} \quad \text{perplexity}(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i \mid w_{i-1})}}$$

$$\log(p_1 \times p_2 \times p_3 \times p_4) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V} \quad P_{\text{Laplace}}(w_n \mid w_{n-1}) = \frac{C(w_{n-1} w_n) + 1}{\sum_w C(w_{n-1} w) + 1} = \frac{C(w_{n-1} w_n) + 1}{C(w_{n-1}) + V}$$

$$P_{\text{Add-k}}^*(w_n \mid w_{n-1}) = \frac{C(w_{n-1} w_n) + k}{C(w_{n-1}) + kV}$$

If we are trying to compute $P(w_n \mid w_{n-2} w_{n-1})$ but we have no examples of a particular trigram $w_{n-2} w_{n-1} w_n$, we can instead estimate its probability by using the bigram probability $P(w_n \mid w_{n-1})$. Similarly, if we don't have counts to compute $P(w_n \mid w_{n-1})$, we can look to the unigram $P(w_n)$.

In other words, sometimes using **less context** is a good thing, helping to generalize more for contexts that the model hasn't learned much about. There are two ways to use this n-gram "hierarchy". In **backoff**, we use the trigram if the evidence is sufficient, otherwise we use the bigram, otherwise the unigram. In other words, we only "back off" to a lower-order n-gram if we have zero evidence for a higher-order n-gram. By contrast, in **interpolation**, we always mix the probability estimates from all the n-gram estimators, weighting and combining the trigram, bigram, and unigram counts.

Language model computes:

Probability of the sentence

$$P(W) \text{ or } P(w_1 w_2 \dots w_{n-1} w_n)$$

$P(\text{"its water is so transparent"})$

$$= P(\text{its}) \times P(\text{water} \mid \text{its}) \times P(\text{is} \mid \text{its water})$$

$$\times P(\text{so} \mid \text{its water is})$$

$$\times P(\text{transparent} \mid \text{its water is so})$$

$$P(w_1 w_2 \dots w_n) = \prod_{k=1}^n P(w_k \mid w_1^{k-1})$$

<s>

I am Sam

</s>

Sam I am

</s>

I do not like green eggs and ham </s>

</s>

do

</s>

I

</s>

am

</s>

the

</s>

green

</s>

eggs

</s>

and

</s>

ham

</s>

and

</s>

the

</s>

ham

</s>

the

Question 2. N-gram language models [5 Marks]

"Erring is human. To err is human. Seeing is believing. To see is to believe."

Given above training data, use bigram language model to compare the likelihood of below test sentences.

Note: To handle zero probability, use interpolation with weights 0.3 & 0.7 for unigram and bigram respectively. Don't account for punctuation ie., "." for the modelling and predictions.

- a) Test Sentence 1: Seeing is human
- b) Test Sentence 2: To believe is human

Cleaned text: "erring is human to err is human seeing is believing to see is to believe"

unigrams

$$P(\text{erring}) = \frac{1}{15}$$

$$P(\text{is}) = \frac{4}{15}$$

$$P(\text{human}) = \frac{2}{15}$$

$$P(\text{to}) = \frac{3}{15} = \frac{1}{5}$$

$$P(\text{err}) = \frac{1}{15}$$

$$P(\text{believing}) = \frac{1}{15}$$

$$P(\text{believe}) = \frac{1}{15}$$

$$P(\text{see}) = \frac{1}{15}$$

$$P(\text{seeing}) = \frac{1}{15} \text{ & } P(<\text{s}>) = P(<\text{E}>) = \frac{1}{15}$$

bigrams

$$P(\text{erring} | <\text{s}>) = \frac{1}{1}$$

$$P(\text{is} | \text{erring}) = \frac{1}{1}$$

$$P(\text{human} | \text{is}) = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{to} | \text{human}) = \frac{1}{2}$$

$$P(\text{err} | \text{to}) = \frac{1}{3}$$

$$P(\text{is} | \text{err}) = \frac{1}{1}$$

$$P(\text{seeing} | \text{human}) = \frac{1}{2}$$

$$P(\text{see} | \text{to}) = \frac{1}{2}$$

$$P(\text{is} | \text{see}) = \frac{1}{1}$$

$$P(\text{to} | \text{is}) = \frac{1}{4}$$

$$P(\text{believe} | \text{to}) = \frac{1}{3}$$

$$P(<\text{E}> | \text{believe}) = \frac{1}{1}$$

Test sentence 1 : Seeing is human

$$P(\text{"seeing is human"}) = \hat{P}(\text{seeing} | <\text{s}>) \times \hat{P}(\text{is} | \text{seeing}) \\ \times \hat{P}(\text{human} | \text{is}) \times \hat{P}(<\text{E}> | \text{human})$$

$$(0.3 \times P(\text{seeing}) + 0.7 \times P(\text{seeing} | <\text{s}>)) \times = (0.3 \times \frac{1}{15} + 0.7 \times 0) \times = 0.00012168$$

$$(0.3 \times P(\text{is}) + 0.7 \times P(\text{is} | \text{seeing})) \times = (0.3 \times \frac{4}{15} + 0.7 \times 1) \times = 1.2168e^{-4}$$

$$(0.3 \times P(\text{human}) + 0.7 \times P(\text{human} | \text{is})) \times = (0.3 \times \frac{2}{15} + 0.7 \times \frac{1}{2}) \times$$

$$(0.3 \times P(<\text{E}>) + 0.7 \times P(<\text{E}> | \text{human})) \times = (0.3 \times \frac{1}{15} + 0.7 \times 0)$$

```
In [11]: (0.3/15)*(0.3*4/15 + 0.7)*(0.3*2/15 + 0.7/2)*(0.3/15)
Out[11]: 0.0001216799999999999
```

Question 4. Vector semantics [1.5+1.5=3 Marks]

Answer the following: with respect to Training Corpus D:

Training Corpus D

document - d1: Today is sunny day! I like sunny days
 document - d2: I do not like sunny days. I like chocolate
 document - d3: I like both sunny and rainy day

Use the TF-IDF to decide on :

The word "sunny" is discriminative in identifying documents in corpus D.
 The word "chocolate" is discriminative in identifying documents in corpus D. Justify your answer.

TF is calculated by dividing the number of times a word appears in the document by the total number of words in that document.

TF	D1	D2	D3
Sunny	2/8=0.25	1/9=0.1111	1/7=0.1429
Chocolate	0	1/9=0.1111	0

IDF is calculated by taking the logarithm of the ratio of the total number of documents (3) to the number of documents that contain word

	IDF
sunny	$\log_{10}(3/3) = 0$
chocolate	$\log_{10}(3/1) = 0.47712125471966244$

TF-IDF is calculated by multiplying the TF and IDF values

TF-IDF	D1	D2	D3
Sunny	0	0	0
Chocolate	0	(1/9)*log ₁₀ (3/1)=0.0530134727466291	0

- a) The word "sunny" is discriminative in identifying documents in corpus D. (No)
- b) The word "chocolate" Is discriminative In identifying documents in corpus D. (Yes)

Since TF-IDF value of "chocolate" in D2 is more than that of "sunny"

```
d1 = "today is sunny day I like sunny days"
d2 = "I do not like sunny days I like chocolate"
d3 = "I like both sunny and rainy day"

corpus = [d1,d2,d3]

import pandas as pd
import numpy as np

words_set = set()

for doc in corpus:
    words = doc.split(' ')
    words_set = words_set.union(set(words))

print('Number of words in the corpus:',len(words_set))
print('The words in the corpus: \n', words_set)

Number of words in the corpus: 13
The words in the corpus:
{'today', 'day', 'both', 'do', 'chocolate', 'not', 'rainy', 'days', 'and', 'is', 'like', 'sunny', 'I'}
```

```
n_docs = len(corpus)      # Number of documents in the corpus
n_words_set = len(words_set) # Number of unique words in the corpus

df_tf = pd.DataFrame(np.zeros((n_docs, n_words_set)), columns=list(words_set))

# Compute Term Frequency (TF)
for i in range(n_docs):
    words = corpus[i].split(' ') # Words in the document
    for w in words:
        df_tf[w][i] = df_tf[w][i] + (1 / len(words))

df_tf
```

	today	day	both	do	chocolate	not	rainy	days	and	is	like	sunny	I
0	0.125	0.125000	0.000000	0.000000	0.000000	0.000000	0.000000	0.125000	0.000000	0.125	0.125000	0.250000	0.125000
1	0.000	0.000000	0.000000	0.111111	0.111111	0.111111	0.000000	0.111111	0.000000	0.000	0.222222	0.111111	0.222222
2	0.000	0.142857	0.142857	0.000000	0.000000	0.000000	0.142857	0.000000	0.142857	0.000	0.142857	0.142857	0.142857

```
print("IDF of: ")
idf = {}

for w in words_set:
    # no.of documents in the corpus that contain this word
    k = 0

    for i in range(n_docs):
        if w in corpus[i].split():
            k += 1

    idf[w] = np.log10(n_docs / k)
    print(w,'is log({}/{})'.format(n_docs,k),f'= {idf[w]}')

df_tf_idf = df_tf.copy()
```

```
for w in words_set:
    for i in range(n_docs):
        df_tf_idf[w][i] = df_tf[w][i] * idf[w]

df_tf_idf
```

	today	day	both	do	chocolate	not	rainy	days	and	is	like	sunny	I
0	0.05964	0.022011	0.00000	0.000000	0.000000	0.000000	0.00000	0.022011	0.00000	0.05964	0.0	0.0	0.0
1	0.00000	0.000000	0.00000	0.053013	0.053013	0.053013	0.00000	0.019566	0.00000	0.00000	0.0	0.0	0.0
2	0.00000	0.025156	0.06816	0.000000	0.000000	0.000000	0.06816	0.000000	0.06816	0.00000	0.0	0.0	0.0

Skip Gram

Given a training corpus: "I live in Bangalore, India", use the skip-gram negative sampling method and answer the following: The initial embedding matrix has dimensions $|V| \times 3$ and is given as follows:

- Generate the training dataset for an input word 'I' and context window of 1 next word and hyper parameter value $k=3$ for the negative sampling task. Use the information available in the question.
- Calculate the error for the above dataset during the first iteration of skip-gram training.

I	0.03	-0.504	0.28
Live	0.1	-0.009	-0.8
In	-0.56	0.745	0.002
Bangalore	-0.592	-0.81	23
India	-0.728	-0.06	0.499

The initial context matrix has $|V| \times 3$ dimension and given as follows:

I	-0.572	-0.588	-0.501
Live	0.22	0.83	-0.72
In	-0.84	0.501	0.246
Bangalore	-0.422	0.911	0.55
India	-0.087	-0.45	-0.07

Negative Sampling I want a glass of orange juice to go along with my cereal.

for a target word

- pick 1 context word and mark it as +ve sample
- pick 'K' context word and mark them as -ve sample even though they are in context window

a) for word "I" given
immediate word is live $K=3$

context	target
live	$t = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
in	
bangalore	
india	

picking remaining words
 $K=3$
as -ve

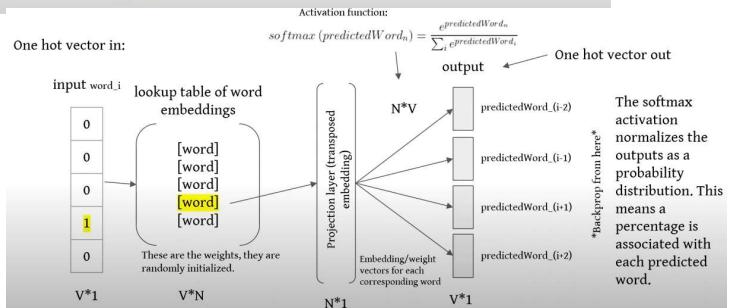
b.) since we got context words as "live", "in", "bangalore", "India" we will pick context vectors of these context words and pick embedding vector of "I" as those are the weights that we want to learn.

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 0.03 & -0.504 & 0.28 \\ 0.1 & -0.009 & -0.8 \\ -0.56 & 0.745 & 0.002 \\ -0.592 & -0.81 & 23 \\ -0.728 & -0.06 & 0.499 \end{bmatrix} \rightarrow \begin{bmatrix} 0.03 & -0.504 & 0.28 \end{bmatrix} \times \begin{bmatrix} 0.22 & 0.83 & -0.72 \\ -0.84 & 0.501 & 0.246 \\ -0.422 & 0.911 & 0.55 \\ -0.087 & -0.45 & -0.07 \end{bmatrix}_{1 \times 3}^{N \times V} \xrightarrow{\text{Transpose}} \begin{bmatrix} 0.22 & -0.84 & -0.422 & -0.087 \\ 0.83 & 0.501 & 0.911 & 0.55 \\ -0.422 & 0.911 & 0.55 & -0.07 \end{bmatrix}_{3 \times N}^{N \times 1}$$

$$\begin{bmatrix} 0.03 & -0.504 & 0.28 \end{bmatrix} \times \begin{bmatrix} 0.22 & -0.84 & -0.422 & -0.087 \\ 0.83 & 0.501 & 0.911 & 0.55 \\ -0.422 & 0.911 & 0.55 & -0.07 \end{bmatrix}_{3 \times N}^{N \times 1} = \begin{bmatrix} 0.61332, -0.208824, -0.317804 \\ 0.03 \times 0.22 + -0.504 \times 0.83 + 0.28 \times -0.72 \\ 0.03 \times -0.84 + -0.504 \times 0.501 + 0.28 \times 0.246 \\ 0.03 \times -0.422 + -0.504 \times 0.911 + 0.28 \times 0.55 \\ 0.03 \times -0.087 + -0.504 \times -0.45 + 0.28 \times -0.07 \end{bmatrix}_{3 \times 1}$$

$$\sigma(x) = \frac{1}{1+e^{-x}} = \begin{bmatrix} 0.351302, 0.497993, 0.62121, 0.55097 \end{bmatrix} = v$$

$$v - t = \text{error}$$



The softmax activation normalizes the outputs as a probability distribution. This means a percentage is associated with each predicted word.

POS tagging [4 Marks]

Tokenize and tag the following sentences using the Brown Corpus table given,

Sentence 1: "he will chair the session"

Sentence 2: "It is a chair"

Sentence 3: "He will race the car."

Sentence 4: "It is a race."

Write the best sequence of tags for sentence 1, sentence 2, sentence 3 and sentence.

Sr.No	Examples	Tag	Description
1	He, she, It	PRP	Personal Pronoun
2	will	MD	modal
3	Chair, race	VB	verb
4	is	VBZ	Verb, 3 rd person singular
5	a, the	DT	determiner
6	chair, session, race, car	NN	Noun, singular or mass
7	.	DOT	DOT

By manual tagging I know

S PRP MO VB PT NN E
"he will chair the session"

S PRP VB PT NN E
"It is a chair"

S PRP MD VB DT NN E
"He will race the car."

S PRP VB DT NN DOT E
"It is a race."

Emission counts

	PRP	MD	VB	VBZ	DT	NN	DOT
he	2						
she							
it	2						
will		2					
chair			1			1	
race			1			1	
is				2			
a					2		
the					2		
session						1	
car						1	
.							2
Total	4	2	2	2	4	4	2

Emission probabilities

	PRP	MD	VB	VBZ	DT	NN	DOT
he	$2/4 = 0.5$						
she							
it	$2/4 = 0.5$						
will		$2/2 = 1$					
chair			$1/2 = 0.5$			$1/4 = 0.25$	
race			$1/2 = 0.5$			$1/4 = 0.25$	
is				$2/2 = 1$			
a					$2/4 = 0.5$		
the					$2/4 = 0.5$		
session						$1/4 = 0.25$	
car						$1/4 = 0.25$	
.							$2/2 = 1$

Transition counts

	PRP	MD	VB	Vbz	DT	NN	DOT	<END>	Total
<START>	4								4
PRP		2		2					4
MD			2						2
VB					2				2
Vbz						2			2
DT							4		4
NN							2	2	4
DOT								2	2

Transition probabilities

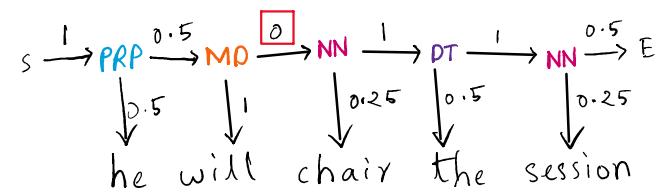
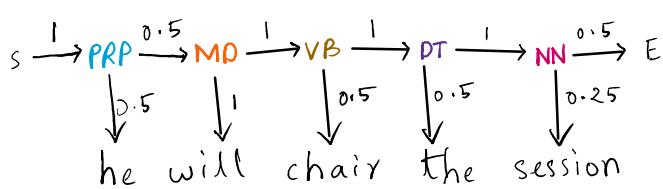
	PRP	MD	VB	Vbz	DT	NN	DOT	<END>
<START>	4/4=1							
PRP		2/4=0.5		2/4=0.5				
MD			2/2=1					
VB					2/2=1			
Vbz						2/2=1		
DT							4/4=1	
NN							2/4=0.5	2/4=0.5
DOT								2/2=1

from the given tag table the words "chair" & "race" have multiple tags and using training sentences we've built Emission & Transition probabilities which inturn can be used in HMM tagger to most find most probable tags

sentence 1 :

PRP MD VB DT NN
he will chair the session

PRP MD NN DT NN
he will chair the session



$$1 \times 0.5 \times 0.5 \times 1 \times 1 \times 0.5 \times 1 \times 0.5 \times 1 \times 0.25 \times 0.5$$

$$= 0$$

$$= 0.0078125$$

since $p(\text{PRP MD VB DT NN}) > p(\text{PRP MD NN DT NN})$ most probable tags are PRP MD VB DT NN
he will chair the session

Sentence 2:

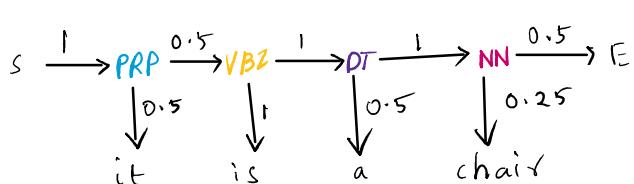
PRP VBZ DT NN

it is a chair

or

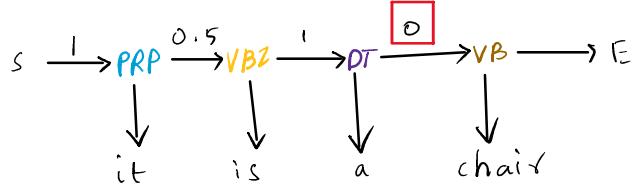
PRP VBZ DT VB

it is a chair



$$= 1 \times 0.5 \times 0.5 \times 1 \times 1 \times 0.5 \times 1 \times 0.25 \times 0.5$$

$$= 0.015625$$



$$= 0$$

since $p(\text{PRP, VBZ, DT, NN}) > p(\text{PRP, VBZ, DT, VB})$ most probable tags are

PRP VBZ DT NN
it is a chair

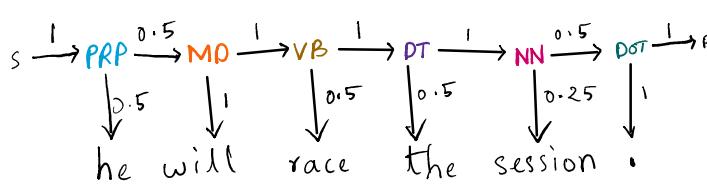
Sentence 3:

PRP MD VB DT NN DOT or

he will race the session .

PRP MD NN DT NN DOT

he will race the session .



$$1 \times 0.5 \times 0.5 \times 1 \times 1 \times 0.5 \times 1 \times 0.5 \times 1 \times 0.25 \times 0.5 \times 1$$

$$= 0$$

$$= 0.0078125$$

since $p(\text{PRP MD VB DT NN}) > p(\text{PRP MD NN DT NN})$ most probable tags are

PRP MD VB DT NN
he will race the session

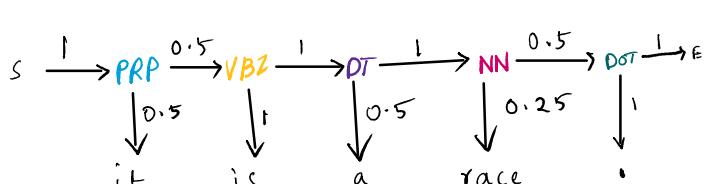
Sentence 4:

PRP VBZ DT NN DOT

it is a race .

PRP VBZ DT VB DOT

it is a race .



$$= 1 \times 0.5 \times 0.5 \times 1 \times 1 \times 0.5 \times 1 \times 0.25 \times 0.5$$

$$= 0$$

$$= 0.015625$$

since $p(\text{PRP, VBZ, DT, NN, DOT}) > p(\text{PRP, VBZ, DT, VB, DOT})$ most probable tags are

PRP VBZ DT NN DOT
it is a race .

```

function VITERBI(observations of len  $T$ ,state-graph of len  $N$ ) returns best-path, path-prob
  create a path probability matrix viterbi[ $N, T$ ]
  for each state  $s$  from 1 to  $N$  do ; initialization step
     $viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$ 
     $backpointer[s, 1] \leftarrow 0$ 
  for each time step  $t$  from 2 to  $T$  do ; recursion step
    for each state  $s$  from 1 to  $N$  do
       $viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$ 
       $backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$ 
     $bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$  ; termination step
     $bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$  ; termination step
  bestpath  $\leftarrow$  the path starting at state bestpathpointer, that follows backpointer[] to states back in time
  return bestpath, bestpathprob

```

transition probabilities $P(t_i | t_{i-1})$

	NNP	MD	VB	JJ	NN	RB	DT
$<s>$	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026
NNP	0.3777	0.0110	0.0009	0.0084	0.0584	0.0090	0.0025
MD	0.0008	0.0002	0.7968	0.0005	0.0008	0.1698	0.0041
VB	0.0322	0.0005	0.0050	0.0837	0.0615	0.0514	0.2231
JJ	0.0366	0.0004	0.0001	0.0733	0.4509	0.0036	0.0036
NN	0.0096	0.0176	0.0014	0.0086	0.1216	0.0177	0.0068
RB	0.0068	0.0102	0.1011	0.1012	0.0120	0.0728	0.0479
DT	0.1147	0.0021	0.0002	0.2157	0.4744	0.0102	0.0017

Observation likelihoods B

	Janet	will	back	the	bill	row sum
NNP	0.000032	0	0	0.000048	0	should be 1
MD	0	0.308431	0	0	0	but here its
VB	0	0.000028	0.000672	0	0.000028	not thus it
JJ	0	0	0.000340	0	0	is only a subset
NN	0	0.000200	0.000223	0	0.002337	
RB	0	0	0.010446	0	0	
DT	0	0	0	0.506099	0	

$\langle s \rangle$

$\overbrace{\begin{array}{l} DT \\ RB \\ NN \\ JJ \\ VB \\ MD \end{array}}^{\text{should be } 1}$

$p_e(\text{Janet} | \langle s \rangle) = 0 \therefore \text{ignore}$

$viterbi[\text{Janet}, \text{NNP}] = p_t(\text{NNP} | \langle s \rangle) \times p_e(\text{Janet} | \text{NNP})$

$= 0.2767 \times 0.000032$

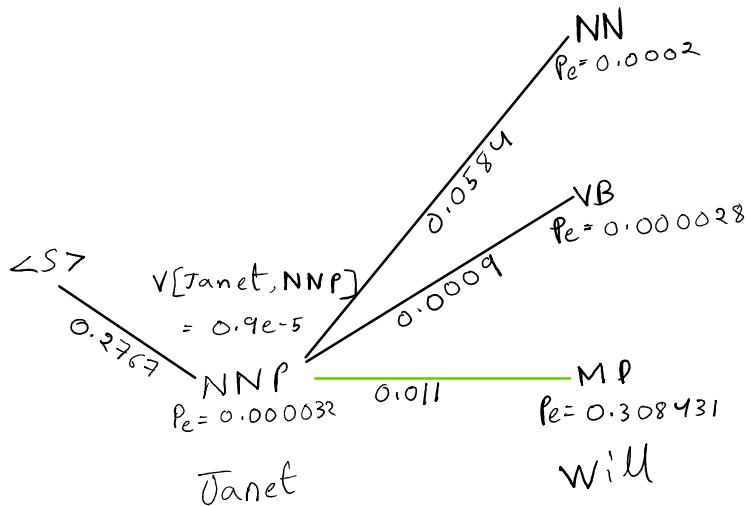
$= 0.0000088544$

$V[\text{Janet}, \text{NNP}] \approx 0.00009 = 0.9e-5$

$$V[will, NN] = V[Janet, NNP] \times P_t(NN|will) \times P_e(will|NN)$$

$$= 0.9e-5 \times 0.0584 \times 0.2e-3$$

$$= 1.0512e-10$$

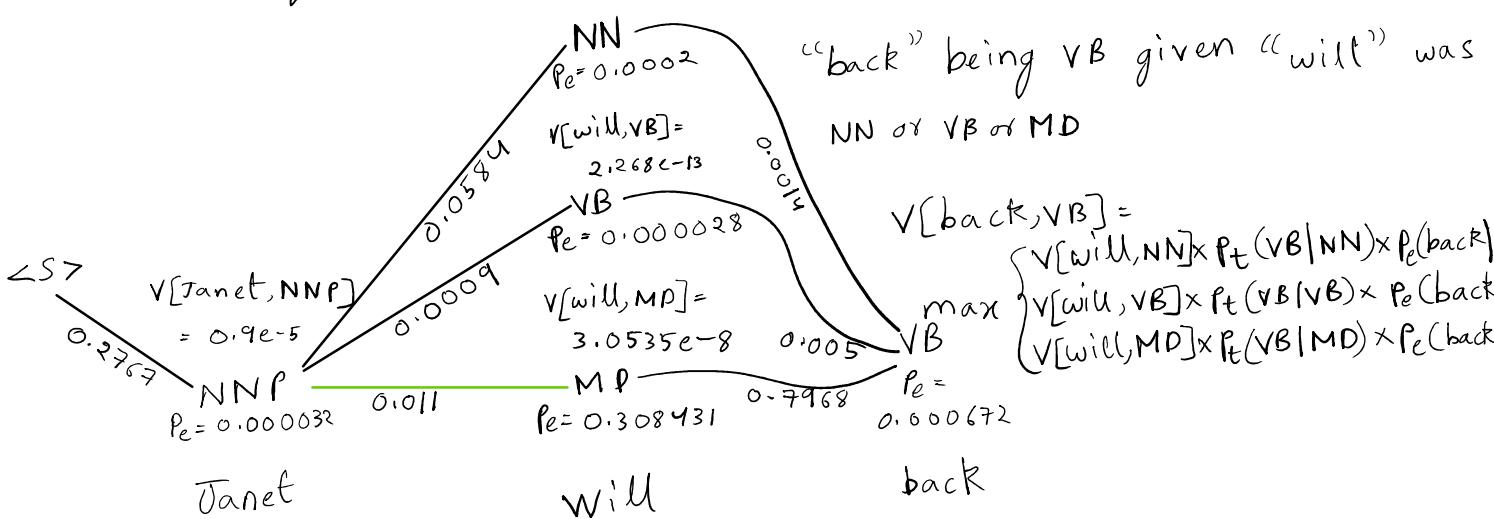


since out of $V[will, NN]$, $V[will, VB]$, $V[will, MD]$
 $V[will, MD]$ is highest. This indicates in given sentence word "will" belongs to MD tag
if previous word ("janet") was NNP tag most probably.

$$V[will, NN] = 0.9e-5 \times 0.011 \times 0.308431$$

$$= 3.0534669e-8$$

$$\approx 3.0535e-8$$



There are 3 possibilities for

"back" being VB given "will" was NN or VB or MD

$$V[back, VB] = \max \{ V[will, NN] \times P_t(VB|NN) \times P_e(back|VB), V[will, VB] \times P_t(VB|VB) \times P_e(back|VB), V[will, MD] \times P_t(VB|MD) \times P_e(back|VB) \}$$

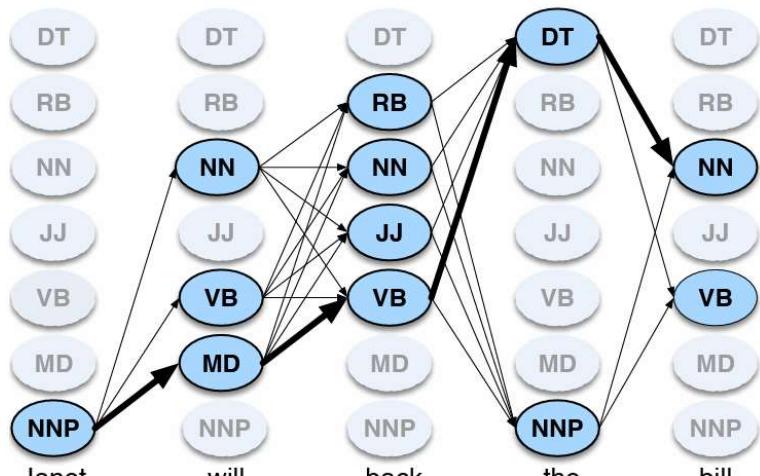
$$= \max \{ 1.0512e-10 \times 0.0014 \times 0.000672, 2.268e-13 \times 0.005 \times 0.000672, 3.0535e-8 \times 0.7968 \times 0.000672 \}$$

$$= \max \{ 9.89e-17, 7.6208e-19, 1.635e-11 \} \quad V[back, VB] = 1.635e-11$$

similarly we can calculate

$V[back, DT]$, $V[back, NN]$, $V[back, RB]$ and other viterbi values

and finally we will get beside network of paths. Key idea is every time we move from one word to another we eliminate some possible paths



Question 7. HMM [4 Marks]

By using Viterbi Algorithm, find whether the word “drive” is VB or NN through computing their probabilities for the sentence, “I love to drive”. The transition probabilities and the observation likelihood for this corpus are as follows:

	VB	TO	NN	PRP
<S>	0.019	0.0043	0.041	0.067
VB	0.0038	0.035	0.047	0.0070
TO	0.83	0	0.0047	0
NN	0.0040	0.016	0.087	0.0045
PRP	0.23	0.00079	0.0012	0.00014

	I	Love	to	drive
VB	0	0.0082	0	.07
TO	0	0	0.88	0
NN	0	0.0074	0	0.057
PRP	0.48	0	0	0

Assume the values for the third column of the Viterbi table which corresponds to observation 3 (word is “to”) and they are VB: 0, TO: 0.018, NN: 0, PPSS: 0.

(Hint: Show the computations of the last column of the Viterbi table).

VB		0	$v[4,4] = v[3,3]*P(VB TO)*P(drive VB)$ $0.018*0.83*0.07 = 0.0010458$	
TO		$0.018 = v[3,3]$	$v[4,3] = v[3,3]*P(TO TO)*P(drive TO)$ $0.018*0*0 = 0$	
NN		0	$v[4,2] = v[3,3]*P(NN TO)*P(drive NN)$ $0.018*0.0047*0.057 = 0.0000048222$	
PRP		0	$v[4,1] = v[3,3]*P(PRP TO)*P(drive PRP)$ $0.018*0*0 = 0$	
	I	Love	To	Drive

$v[4,4] = P(\text{drive} = \text{VB}) = 0.0010458$ and $v[4,2] = P(\text{drive} = \text{NN}) = 0.0000048222$
Since $P(\text{drive} = \text{VB}) > P(\text{drive} = \text{NN})$ the word drive is VB in the given sentence

Top Down Parser

1 The 2 old 3 man 4 cried 5

1. $S \rightarrow NP\ VP$
2. $NP \rightarrow ART\ N$
3. $NP \rightarrow ART\ ADJ\ N$

4. $VP \rightarrow V$
5. $VP \rightarrow V\ NP$

the: ART
old: ADJ, N
man: N, V
cried: V

Step	Current State	Backup States	Comment
1.	((S) 1)		
2.	((NP VP) 1)		S rewritten to NP VP
3.	((ART N VP) 1)	((ART ADJ N VP) 1)	NP rewritten producing two new states
4.	((N VP) 2)	((ART ADJ N VP) 1)	
5.	((VP) 3)	((ART ADJ N VP) 1)	the backup state remains
6.	((V) 3)	((V NP) 3) ((ART ADJ N VP) 1)	
7.	(() 4)	((V NP) 3) ((ART ADJ N VP) 1)	
8.	((V NP) 3)	((ART ADJ N VP) 1)	the first backup is chosen
9.	((NP) 4)	((ART ADJ N VP) 1)	
10.	((ART N) 4)	((ART ADJ N) 4) ((ART ADJ N VP) 1)	looking for ART at 4 fails
11.	((ART ADJ N) 4)	((ART ADJ N VP) 1)	fails again
12.	((ART ADJ N VP) 1)		now exploring backup state saved in step 3
13.	((ADJ N VP) 2)		
14.	((N VP) 3)		
15.	((VP) 4)		
16.	((V) 4)	((V NP) 4)	
17.	(() 5)		success!

The algorithm starts with the initial state $((S) 1)$ and no backup states.

1. Select the current state: Take the first state off the possibilities list and call it C. If the possibilities list is empty, then the algorithm fails (that is, no successful parse is possible).
2. If C consists of an empty symbol list and the word position is at the end of the sentence, then the algorithm succeeds.
3. Otherwise, generate the next possible states.
 - 3.1. If the first symbol on the symbol list of C is a lexical symbol, and the next word in the sentence can be in that class, then create a new state by removing the first symbol from the symbol list and updating the word position, and add it to the possibilities list.
 - 3.2. Otherwise, if the first symbol on the symbol list of C is a non-terminal, generate a new state for each rule in the grammar that can rewrite that nonterminal symbol and add them all to the possibilities list.

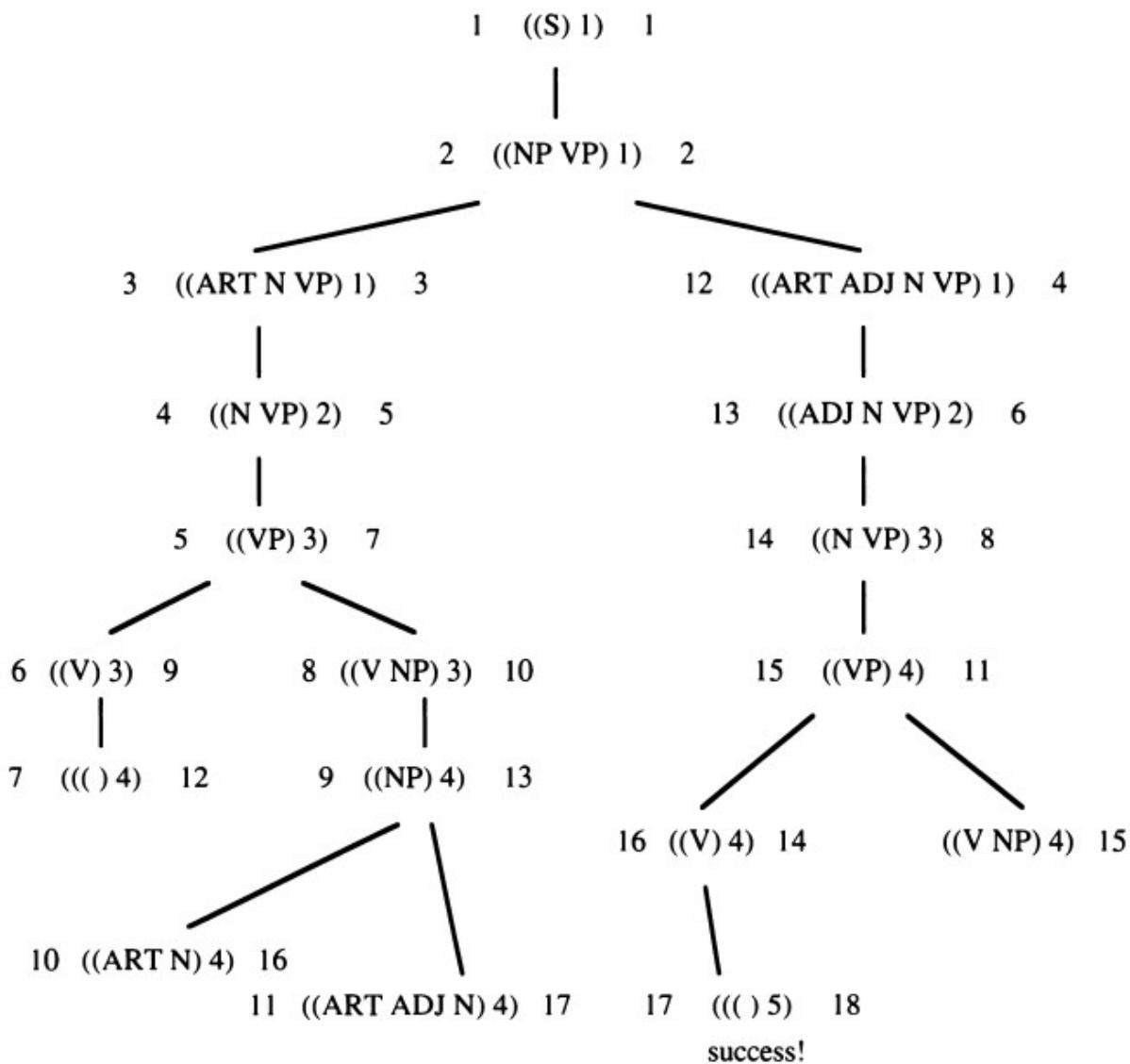


Figure 3.7 Search tree for two parse strategies (depth-first strategy on left; breadth-first on right)

Top Down Parser (More examples)

People laugh
1 2 3 These are positions

Lexicon:

People - N, V
Laugh - N, V

This indicate that both Noun and Verb is possible for the word “People”

- Grammar – $S \rightarrow NP\ VP$
 $NP \rightarrow DT\ N\ | N$
 $VP \rightarrow V\ ADV\ | V$

State	Backup State	Action
1. $((S) 1)$ search for 1st phrase	-	-
2. $((NP VP) 1)$ word at index 1 i.e "people" is it Noun phrase(NP)?	by given Lexicon it's a Noun(N)	-
3a. $((DT N VP) 1)$ "people" $\in \{N, V\}$ $DT \& \dots$ next state	$((N VP) 1)$	-
3b. $((N VP) 1)$ search phrase = N & current index word "people" $\in \{N, V\}$	Match found	Verb(V)
4. $((VP) 2)$ remove matched search phrase & index++	Consume "People"	-
5a. $((V ADV) 2)$ search phrase V is in $((V) 2)$	-	-
6. $((ADV) 3)$ current "laugh" $\in \{N, V\}$ remove search phrase V & index++	Consume "laugh"	-
5b. $((V) 2)$ at index 3 reached end i.e empty string & it doesn't have any productions or lexicon match so backup	Consume "laugh"	-
6. $((.) 3)$ state is picked backup state moved to index 2 "laugh" & search phrase(sp) V	Consume "laugh"	-

"₁ The ₂ old ₃ man ₄ cried ₅"

1. $S \rightarrow NP\ VP$	4. $VP \rightarrow V$
2. $NP \rightarrow ART\ N$	5. $VP \rightarrow V\ NP$
3. $NP \rightarrow ART\ ADJ\ N$	

cried: V
old: ADJ, N
man: N, V
the: ART

Step	Current State	Backup States	Comment	Step	Current State	Backup States	Comment
1.	((S) 1)		initial position	9.	((NP) 4)		
2.	((NP VP) 1)		S rewritten to NP VP	10.	((ART N) 4)	((ART ADJ N VP) 1)	
3.	((ART N VP) 1)		NP rewritten by rules 2&3	11.	((ART ADJ N) 4)	((ART ADJ N) 4)	looking for ART fails
4.	((N VP) 2)	((ART ADJ N VP) 1)		12.	((ART ADJ N VP) 1)	((ART ADJ N VP) 1)	
5.	((VP) 3)	((ART ADJ N VP) 1)		13.	((ADJ N VP) 2)		fails again
6.	((V) 3)	((ART ADJ N VP)	VP rewritten by rules 4&5	14.	((N VP) 3)		
		1) ((V NP) 3)		15.	((VP) 4)		
		((ART ADJ N VP) 1)		16.	((V) 4)		
7.	(() 4)						
				17.	(())		
8.	((V NP) 3)		the first backup is chosen				
				5)			success!

Top Down Parser

Ref https://gawron.sdsu.edu/compling/course_core/lectures/simple_top_down.htm

Think of this as an algorithm that takes one particular view of how to do a derivation. Suppose we are trying to parse the sentence *John likes Richard M. Nixon* with respect to the grammar:

- S -> NP VP | VP
- VP -> V NP | V PP
- NP -> Art (Adj) N | PN
- VP -> is/are Adj | V NP
- PP -> P NP
- PN -> Richard M Nixon | John
- V -> likes | runs | eats
- P -> in | on | under
- Art -> the | a
- Adj -> red | big
- N -> house | garden | car

Non-terminals are symbols that occur on the LHS of a rule. The non-terminals in this grammar

S, NP, VP, PP, P, Art, Adj, N, V, PN

Terminals are symbols that occur ONLY on the RHS of a rule. The non-terminals in this grammar

is are Richard M Nixon John likes runs
eats in on under the a red big

The **pre-terminals** are symbols that occur on the LHS of a rule introducing ONLY terminals. The pre-terminals are:

N Art Adj PN V P

Step	Curr. State	State Stack	Rewrite/Found
1.	((S) (John likes R.M.N.))		initial state
2.	((NP VP) (John likes R.M.N.))	((VP) (John likes R.M.N.))	S -> NP VP (current state) S -> VP (state stack)
3.	((Art N VP) (John likes R.M.N.))	((Art Adj N VP) (John likes R.M.N.)) ((PN VP) (John likes R.M.N.)) ((VP) (John likes R.M.N.))	NP -> Art N (current state) NP -> Art Adj N (state stack) NP -> PN (state stack)
4.	((Art Adj N VP) (John likes R.M.N.))	((PN VP) (John likes R.M.N.)) ((VP) (John likes R.M.N.))	Popped state stack!! "John" is not an Art
5.	((PN VP) (John likes R.M.N.))	((VP) (John likes R.M.N.))	Popped state stack!! "John" is not an Art
6.	((VP) (likes R.M.N.))	((VP) (John likes R.M.N.))	Lexical pop! "John" is a PN
7.	((V NP) (likes R.M.N.))	((V PP) (likes R.M.N.)) ((VP) (John likes R.M.N.))	VP -> V NP VP -> V PP
8.	((NP) (R.M.N.))	((V PP) (likes R.M.N.)) ((VP) (John likes R.M.N.))	Lexical pop! "likes" is a V
9.	((Art N) (R.M.N.))	((Art Adj N) (R.M.N.)) ((PN) (R.M.N.)) ((V PP) (likes R.M.N.)) ((VP) (John likes R.M.N.))	NP -> Art N (current state) NP -> Art Adj N (backup state) NP -> PN (backup state)
10.	((Art Adj N) (R.M.N.))	((PN) (R.M.N.)) ((V PP) (likes R.M.N.)) ((VP) (John likes R.M.N.))	Popped state stack!! "Richard" is not an Art
11.	((PN) (R.M.N.))	((V PP) (likes R.M.N.)) ((VP) (John likes R.M.N.))	Popped state stack!! "Richard" is not an Art
12.	(O O)	((V PP) (likes R.M.N.)) ((VP) (John likes R.M.N.))	Lexical pop! "R.M.N." is a PN Done!!

Bottom up Chart Parser

For rule $NP \rightarrow ART\ ADJ\ N$ in a bottom-up parser you use the rule to take a sequence $ART\ ADJ\ N$ that you have found and identify it as an NP . A simple bottom-up parser would be expensive, as the parser would tend to try the same matches again and again, thus duplicating much of its work unnecessarily.

To avoid this problem, a data structure called a **chart** is introduced that allows the parser to store the partial results of the matching it has done so far so that the work need not be reduplicated.

Sentence:

"The Large can can hold the water"

Lexicons:

The \rightarrow ART

Large \rightarrow ADJ

Can \rightarrow N, V, AUX

Hold \rightarrow N, V

Water \rightarrow N, V

Grammar Rule (G_Rule):

1. $S \rightarrow NP\ VP$
2. $NP \rightarrow ART\ ADJ\ N$
3. $NP \rightarrow ART\ N$
4. $NP \rightarrow ADJ\ N$
5. $VP \rightarrow AUX\ VP$
6. $VP \rightarrow V\ NP$

INITIALIZE:
set Agenda = list of all possible categories of each input word
(in order of input);
set n = length of input;
set Chart = ();

ITERATE:

loop
if $Agenda = ()$ **then**
 if there is at least one S constituent from 0 to n **then**
 return SUCCESS \circlearrowleft
 else
 return FAIL \circlearrowright
 end if
else ...

Set $C_{ij} = \text{First}(Agenda)$; /* Remove first item from agenda. */
/* C_{ij} is a completed constituent of type C from position i to position j */
Add C_{ij} to **Chart**;

ARC UPDATE:

- a. BOTTOM-UP ARC ADDITION (PREDICTION):
for each grammar rule $X \rightarrow C\ X_1 \dots X_N$ **do**
 Add arc $X \rightarrow C \cdot X_1 \dots X_N$, from i to j , to **Chart**;
 - b. ARC EXTENSION (FUNDAMENTAL RULE):
for each arc $X \rightarrow X_1 \dots \bullet C \dots X_N$, from k to i , **do**
 Add arc $X \rightarrow X_1 \dots C \cdot \dots X_N$, from k to j , to **Chart**;
 - c. ARC COMPLETION:
for each arc $X \rightarrow X_1 \dots X_N C \bullet$ added in step (a) or step (b) **do**
 Move completed constituent X to **Agenda**;
- end if**
end loop

S1 (rule 1 with NP1 and VP2)

S2 (rule 1 with NP2 and VP2)

VP3 (rule 5 with AUX1 and VP2)

NP2 (rule 4) *Added at ③*

VP2 (rule 5)

NP1 (rule 2) *Added at ③*

VP1 (rule 6)

N1

N2

NP3 (rule 3)

V1

V2

V3

V4

ART1	ADJ1	AUX1	AUX2	N3	ART2	N4
-------------	-------------	-------------	-------------	-----------	-------------	-----------

1 the 2 large 3 can 4 can 5 hold 6 the 7 water

start parsing
from 1

① "the" is a ART pick grammar production starting with ART i.e rule 2 & 3 APP new ARCs

A1: $NP \rightarrow ART \cdot ADJ \cdot N \xrightarrow{2} \cdot$ indicates that lexicons to left are seen/encountered and

the ones to right are yet to be seen.

A2: $NP \rightarrow ART \cdot N \xrightarrow{2} \cdot$ - No ARCs can be joined & No ARCs have completed so continue loop

② "large" is an ADJ, pick production starting with ADJ i.e rule 4 & APP a new ARC

A2: $NP \rightarrow ADJ \cdot N \xrightarrow{2} \cdot$ In this ARC an ADJ is seen and one of previous

from A1 $\xleftarrow{A1}$ ARC is searching for an ADJ so these two ARCs can be joined

A3: $A1 + A3 \xrightarrow{1} NP \rightarrow ART\ ADJ \cdot N \xrightarrow{3} \cdot$ No ARC completed so continue loop

③ "can" is an N, V, AUX so production rules 5 & 6 are picked & APP new ARCs

Also previous ARCs A3 & A4 also looking for lexicon N till current index 3 since "can" is N

new Au: $NP \rightarrow ART\ ADJ\ N \xrightarrow{3} \cdot$

new A3: $NP \rightarrow ADJ\ N \xrightarrow{0} \cdot$

since 0 pointer reached end means entire production is seen & its LHS can be added to agenda for reprocessing in index 0

S1 (rule 1 with NP1 and VP2)

S2 (rule 1 with NP2 and VP2)

VP3 (rule 5 with AUX1 and VP2)

NP2 (rule 4) *Added at ③*

VP2 (rule 5)

NP1 (rule 2) *Added at ③*

VP1 (rule 6)

N1

N2

V1

V2

V3

V4

ART1

ADJ1

AUX1

AUX2

N3

ART2

N4

1 the 2 large 3 can 4 can 5 hold 6 the 7 water

$$A1: \frac{NP \rightarrow ART \circ ADJ \circ N}{1 \xrightarrow{1} \xrightarrow{2} 2}$$

$$A2: \frac{NP \rightarrow ART \circ N}{1 \xrightarrow{1} 2}$$

$$A3: \frac{NP \rightarrow ADJ \circ N}{2 \xrightarrow{2} \xrightarrow{3} 3}$$

$$Au: A1 + A3 \quad \frac{NP \rightarrow ART \circ ADJ \circ N}{1 \xrightarrow{1} \xrightarrow{2} \xrightarrow{3} 3}$$

new Au: $\frac{NP \rightarrow ART \circ ADJ \circ N}{1 \xrightarrow{1} \xrightarrow{2} \xrightarrow{3} 3}$ o reached end: add LHS to agenda

new A3: $\frac{NP \rightarrow ADJ \circ N}{2 \xrightarrow{2} \xrightarrow{3} 3}$ o reached end: add LHS to agenda

continue ③ $NP_1 + NP_2 \{NP \text{ at core}\}$

since new lexicons have been added to agenda we have to look for productions starting with NP

$$A5: \frac{S \rightarrow NP \circ VP}{1 \xrightarrow{1} \xrightarrow{4} 4}$$

$$A6: \frac{S \rightarrow NP \circ VP}{2 \xrightarrow{2} \xrightarrow{1} 1}$$

also add ARC's for original lexicons for "can"

$$A7: \frac{VP \rightarrow AUX \circ VP}{3 \xrightarrow{3} \xrightarrow{4} 4} \quad A8: \frac{VP \rightarrow V \circ NP}{3 \xrightarrow{3} \xrightarrow{4} 4}$$

these ARC's didn't reach end
nor can be joined so continue

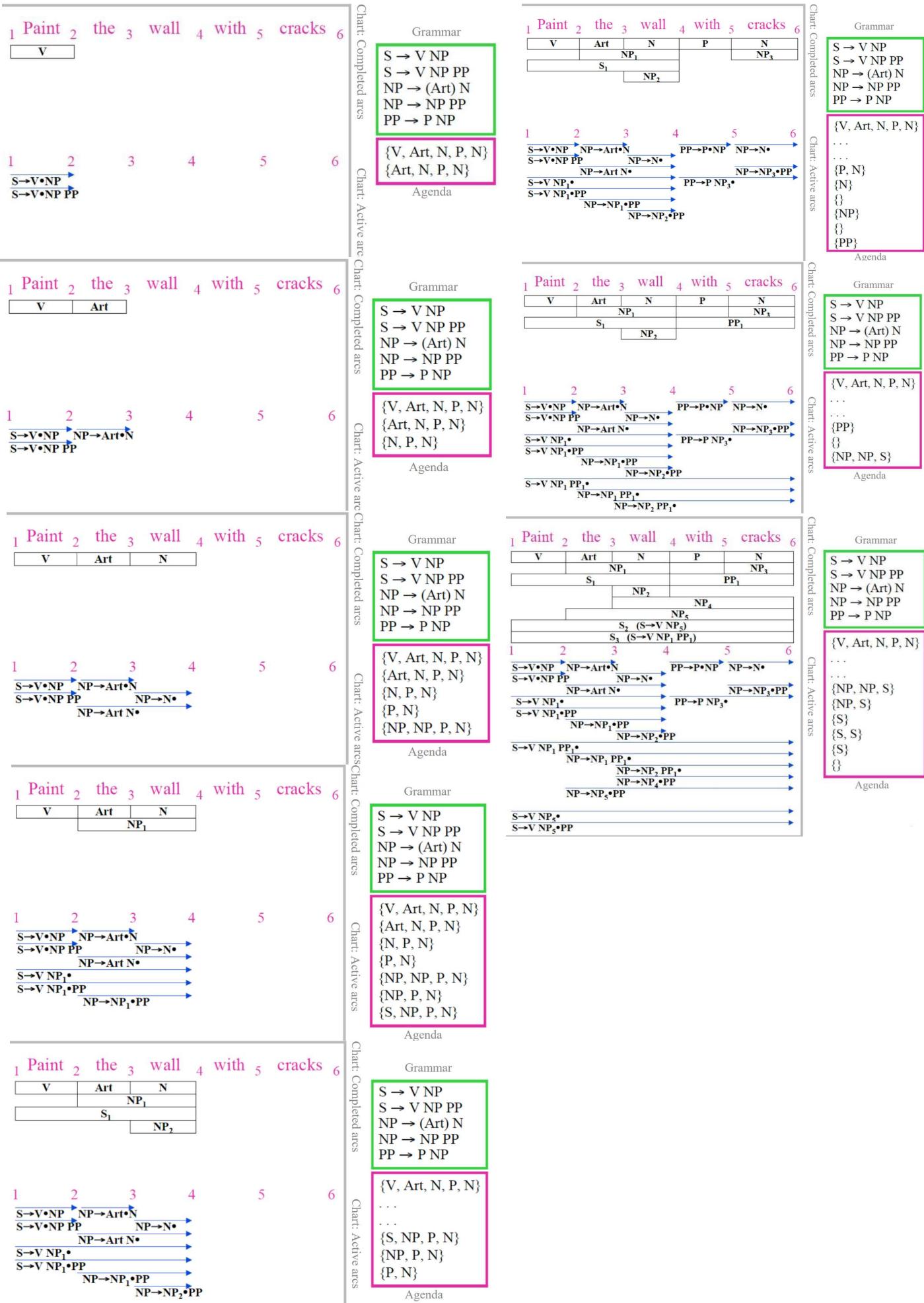
④ ... continue ARC Addition

ARC extension

& ARC completion

until starting non-terminal 'S' is added to agenda spanning from startindex to endindex

Bottom up Chart Parser (more examples)

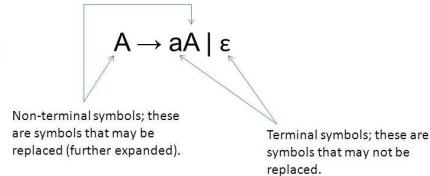


Top Down Chart Parsing/ Earley Parsing

```

function EARLEY-PARSE(words,grammar) returns chart
    ENQUEUE(( $\gamma \rightarrow \bullet S$ , [0,0]),chart[0])
    for i ← from 0 to LENGTH(words) do
        for each state in chart[i] do
            if INCOMPLETE?(state) and NEXT-CAT(state) is not POS then
                PREDICTOR(state)
            elseif INCOMPLETE?(state) and NEXT-CAT(state) is POS then
                SCANNER(state)
            else
                COMPLETER(state)
        end
    end
    return(chart)

```



$\text{R}_7 \quad \text{Det} \rightarrow \text{that} \mid \text{this} \mid \text{a} \mid \text{the}$
 $\text{R}_8 \quad \text{Noun} \rightarrow \text{book} \mid \text{flight} \mid \text{meal} \mid \text{money}$
 $\text{R}_9 \quad \text{Verb} \rightarrow \text{book} \mid \text{include} \mid \text{prefer}$

Rule (R)

$1 \quad S \rightarrow NP \ VP$
 $2 \quad S \rightarrow VP$
 $3 \quad NP \rightarrow \text{Det Nominal}$
 $4 \quad \text{Nominal} \rightarrow \text{Noun}$
 $5 \quad VP \rightarrow \text{Verb}$
 $6 \quad VP \rightarrow \text{Verb NP}$

• Book that flight.

Chart	State	Rule	Start, End	Added By
0	S0	$\gamma \rightarrow \bullet S$	0, 0	Start State
0	S1	$\overline{-S \rightarrow \bullet NP \ VP}$	0, 0	Predictor R1
0	S2	$\overline{-S \rightarrow \bullet VP}$	0, 0	Predictor R2
0	S3	$\overline{NP \rightarrow \bullet \text{Det Nominal}}$	0, 0	Predictor R3
0	S4	$\overline{VP \rightarrow \bullet \text{Verb}}$	0, 0	Predictor R5
0	S5	$\overline{VP \rightarrow \bullet \text{Verb NP}}$	0, 0	Predictor R6

Book • that flight.

if pos right of • then get all rules starting with that pos

1	S6	Verb → book • this also	0, 1	Scanner
1	S7	VP → Verb • move •	0, 1	Completer
1	S8	VP → Verb • NP in parents	0, 1	Completer
1	S9	S → VP • as well	0, 1	Completer
1	S10	NP → • Det Nominal	1, 1	Predictor

Book that • flight.

Since "book" is not pos tag but word so scanner moves • one position ahead

2	S11	Det → that •	1, 2	Scanner
2	S12	NP → Det • Nominal	1, 2	Completer
2	S13	Nominal → • Noun	2, 2	Predictor
3	S14	Noun → flight •	2, 3	Scanner
3	S15	Nominal → Noun •	2, 3	Completer
3	S16	NP → Det Nominal •	1, 3	Completer
3	S17	VP → Verb NP •	0, 3	Completer
3	S18	S → VP •	0, 3	Completer

keep repeating above steps until • hasn't reached end of sentence

Top Down Chart Parsing/ Earley Parsing

Top-down initialization: When category S is the root category of the grammar, then for every rule with S on the LHS, create an empty edge at position [0:0] with a dot leftmost on the RHS. That is, whenever $S \rightarrow A B$, for some A, add an edge [0:0] $S \rightarrow * A B$

Top-down prediction rule: When category X is expected (i.e., just to the right of a dot) then for all rules where X is on the LHS, add a new edge where the dot is leftmost on the RHS with a span of length zero at the current end point. (This rule will be repeated multiple times until no new predictions can be made). That is, whenever $([i:k] A \rightarrow W^* X Y)$ and $X \rightarrow S T$, for some X, add a new edge $[k:k] X \rightarrow * S T$. (Note: W can be empty.)

Fundamental rule: When a new category is seen or nonterminal created, then for all active rules where the category is leftmost on the RHS of the dot, and the spans are adjacent, create a new edge with the dot moved to the left of that category and combine the spans. That is, whenever $([i:j], A \rightarrow X^* Y Z$ and $[j:k] Y \rightarrow S T^*$, for some j and Y, add a new edge $[i:k], A \rightarrow X Y^* Z$ (Note: X and Z can be empty.)

0	$S \rightarrow * NP VBD$ $S \rightarrow * NP VP$	the			
1		dog			
2			likes		
3				meat	
4		0 1 2 3 4			
0	$S \rightarrow * NP VBD$ $S \rightarrow * NP VP$ $NP \rightarrow * DT NN$	the			
1		dog			
2			likes		
3				meat	
4		0 1 2 3 4			
0	$S \rightarrow * NP VBD$ $S \rightarrow * NP VP$ $NP \rightarrow * DT NN$ $DT \rightarrow * the$	the			
1		dog			
2			likes		
3				meat	
4		0 1 2 3 4			
0	$S \rightarrow * NP VBD$ $S \rightarrow * NP VP$ $NP \rightarrow * DT NN$ $DT \rightarrow * the$	the			
1	NN -> * dog	NN -> dog *			
2		VP -> * VBZ NN	likes		
3				meat	
4		0 1 2 3 4			
0	$S \rightarrow * NP VBD$ $S \rightarrow * NP VP$ $NP \rightarrow * DT NN$	NP -> DT * NN	$S \rightarrow NP * VBD$ $S \rightarrow NP * VP$ $NP \rightarrow DT NN *$		
1	NN -> * dog	NN -> dog *			
2		VP -> * VBZ	VBZ -> * likes	VP -> VBZ * N	VBZ -> likes *
3				likes	
4		0 1 2 3 4			
0	$S \rightarrow * NP VBD$ $S \rightarrow * NP VP$ $NP \rightarrow * DT NN$	The	$S \rightarrow NP * VBD$ $S \rightarrow NP * VP$ $NP \rightarrow DT NN *$	$S \rightarrow NP VP *$	
1	NN -> * dog	NN -> dog *			
2		VP -> * VBZ	VBZ -> * likes	VP -> VBZ *	VP -> VBZ NN *
3				likes	
4		0 1 2 3 4			

$S \rightarrow NP VBD$

$S \rightarrow NP VP$

$NP \rightarrow DT NN$

$VP \rightarrow VBZ NN$

$NN \rightarrow dog | meat$

$VBD \rightarrow barked$

$VBZ \rightarrow likes$

$DT \rightarrow the$

“the dog likes meat”.

For each of the sentence rules, make a top-down prediction to create an empty, active edge.

[0:0] $S \rightarrow * NP VBD$

[0:0] $S \rightarrow * NP VP$

[0:0] $NP \rightarrow * DT NN$

[0:1] $DT \rightarrow the ^$

[0:1] $NP \rightarrow DT ^ * NN$

[1:1] $NN \rightarrow * dog$

[1:2] $NN \rightarrow dog ^$

[0:2] $NP \rightarrow DT NN ^ *$

[0:2] $S \rightarrow NP * VBD$

[0:2] $S \rightarrow NP * VP$

[2:2] $VP \rightarrow * VBZ NN$

[2:2] $VBZ \rightarrow * likes$

[2:3] $VBZ \rightarrow likes ^ *$

[2:3] $VP \rightarrow VBZ ^ * NN$

[3:3] $NN \rightarrow * meat$

[3:4] $NN \rightarrow meat ^ *$

[2:4] $VP \rightarrow VBZ NN ^ *$

[0:4] $S \rightarrow NP VP ^ *$

Make more top-down predictions to create active edges for each of the nonterminal categories just to the right of the dot.

Predict *the* and use the fundamental rule to create new edges where the dot in the DT rule and in NP rule move to the right.

Predict *dog*; apply fundamental rule and then make a top down prediction for a VP (using the second S rule.)

Predict *likes*; apply the fundamental rule to create VBZ-> likes *, and again, to add VP -> VBZ * NN

Predict *meat* apply the fundamental rule for meat as a noun (NN) and again to extend the active VP edge, to get VP -> VBZ NN *.

Finally, use the fundamental rule to extend the active S edge, which is now complete.

CFG to CNF

Requires input grammar based on Chomsky Normal Form

– A CNF grammar is a Context-Free Grammar in which:

- Every rule LHS is a non-terminal
- Every rule RHS consists of either a single terminal or two non-terminals.
- Examples:
 - $A \rightarrow B C$
 - $NP \rightarrow N PP$
 - $A \rightarrow a$
 - $Noun \rightarrow man$
- But not:
 - $NP \rightarrow the N$
 - $S \rightarrow VP$

BITS Pilani Deemed to be University

1. Rules that mix terminals and non-terminals on the RHS

– E.g. $NP \rightarrow the Nominal$

– Solution:

• Introduce a dummy non-terminal to cover the original terminal

– E.g. $Det \rightarrow the$

• Re-write the original rule:

– $NP \rightarrow Det Nominal$

3. Rules which have more than two items on the RHS

– E.g. $NP \rightarrow Det Noun PP$

Solution:

– Introduce new non-terminals to spread the sequence on the RHS over more than 1 rule.

- $Nominal \rightarrow Noun PP$
- $NP \rightarrow Det Nominal$

Example:

Check if grammar is in CNF or not, if not convert
 $S \rightarrow NP VP$

$NP \rightarrow ART ADJ N$ (Rule 3 more than 2 non terminal)

$NP \rightarrow a ART N$ (Rule 1 mix of term & non term)

$NP \rightarrow ADJ N$

$P \rightarrow VP$ (Rule 2 unit productions)

$VP \rightarrow AUX VP$

$VP \rightarrow V NP$

(Rule 3) let { let DUM1 $\rightarrow ART ADJ$ } then

$NP \rightarrow DUM1 N$

(Rule 1) let { let DUM2 $\rightarrow a$ } then

$NP \rightarrow DUM2 ART N$

(Rule 3) $NP \rightarrow DUM3 N$,

$DUM3 \rightarrow DUM2 ART$

(Rule 2) { $P \rightarrow VP \rightarrow AUX VP == P \rightarrow AUX VP$ }

{ $P \rightarrow AUX VP$ }

{ $P \rightarrow V NP$ }

To convert a CFG to CNF, we need to deal with three issues:

1. Rules that mix terminals and non-terminals on the RHS
 - E.g. $NP \rightarrow the Nominal$
2. Rules with a single non-terminal on the RHS (called unit productions)
 - E.g. $NP \rightarrow Nominal$
3. Rules which have more than two items on the RHS
 - E.g. $NP \rightarrow Det Noun PP$

8

2. Rules with a single non-terminal on the RHS (called unit productions)

– E.g. $NP \rightarrow Nominal$

– Solution:

• Find all rules that have the form $Nominal \rightarrow ...$

– $Nominal \rightarrow Noun PP$

– $Nominal \rightarrow Det Noun$

– Re-write the above rule several times to eliminate the intermediate non-terminal:

– $NP \rightarrow Noun PP$

– $NP \rightarrow Det Noun$

CNF Grammar

Innovate Achieve

If we parse a sentence with a CNF grammar, we know that:

- Every phrase-level non-terminal (above the part-of-speech level) will have exactly 2 daughters.
 - $NP \rightarrow Det N$
- Every part-of-speech level non-terminal will have exactly 1 daughter, and that daughter is a terminal:

• $N \rightarrow lady$

$S \rightarrow NP VP$

$DUM1 \rightarrow ART ADJ$

$NP \rightarrow DUM1 N$

$DUM2 \rightarrow a$

$DUM3 \rightarrow DUM2 ART$

$NP \rightarrow DUM3 N$

$NP \rightarrow ADJ N$

$P \rightarrow AUX VP$

$P \rightarrow V NP$

CYK Algorithm

CKY: lexical step (j = 1)

	1	2	3	4	5
0	Det				
1					
2					
3					
4					
5					

Lexical lookup

- Matches Det → the

The flight includes a meal.

CKY: lexical step (j = 2)

	1	2	3	4	5
0	Det				
1		N			
2					
3					
4					
5					

Lexical lookup

- Matches N → flight

The flight includes a meal.

CKY: syntactic step (j = 2)

	1	2	3	4	5
0	Det	NP			
1		N			
2					
3					
4					
5					

Syntactic lookup:

- look backwards and see if there is any rule that will cover what we've done so far.

The flight includes a meal.

CKY: lexical step (j = 3)

	1	2	3	4	5
0	Det	NP			
1		N			
2			V		
3					
4					
5					

Lexical lookup

- Matches V → includes

The flight includes a meal.

function CKY-Parse(words, grammar) returns table for

$j \leftarrow \text{from 1 to LENGTH(words)}$ do

loop over the columns

$\text{table}[j-1, j] \leftarrow \{A \mid A \rightarrow \text{words}[j] \in \text{grammar}\}$ fill bottom cell

for $i \leftarrow \text{from } j-2 \text{ downto 0}$ do

fill row i in column j

for $k \leftarrow i+1 \text{ to } j-1$ do

loop over split locations

$\text{table}[i, j] \leftarrow \text{table}[i, j] \cup$ between i and j

$S \rightarrow NP VP \quad \{A \mid A \rightarrow BC \in \text{grammar},$

$B \in \text{table}[i, k]$

$C \in \text{table}[k, j]\}$

Check the grammar for rules that link the constituent in $[i, k]$ with those in $[k, j]$. For each rule found store LHS in cell $[i, j]$.

Det → the

Det → a

N → meal

N → flight

CKY: lexical step (j = 3)

	1	2	3	4	5
0	Det	NP			
1		N			
2			V		
3					
4					
5					

Syntactic lookup

- There are no rules in our grammar that will cover Det, NP, V

The flight includes a meal.

CKY: lexical step (j = 4)

	1	2	3	4	5
0	Det	NP			
1		N			
2			V		
3				Det	
4					
5					

Lexical lookup

- Matches Det → a

The flight includes a meal.

CKY: lexical step (j = 5)

	1	2	3	4	5
0	Det	NP			
1		N			
2			V		
3				Det	
4					N

Lexical lookup

- Matches N → meal

The flight includes a meal.

CYK (contd.)

CKY: syntactic step (j = 5)

	1	2	3	4	5
0	Det	NP			
1		N			
2			V		
3				Det	NP
4					N

Syntactic lookup

- We find that we have NP → Det N

The flight includes a meal.

CKY: syntactic step (j = 5)

	1	2	3	4	5
0	Det	NP			
1		N			
2			V		VP
3				Det	NP
4					N

Syntactic lookup

- We find that we have VP → V NP

The flight includes a meal.

CKY: syntactic step (j = 5)

	1	2	3	4	5
0	Det	NP			S
1		N			
2			V		VP
3				Det	NP
4					N

Syntactic lookup

- We find that we have S → NP VP

The flight includes a meal.

PCFG

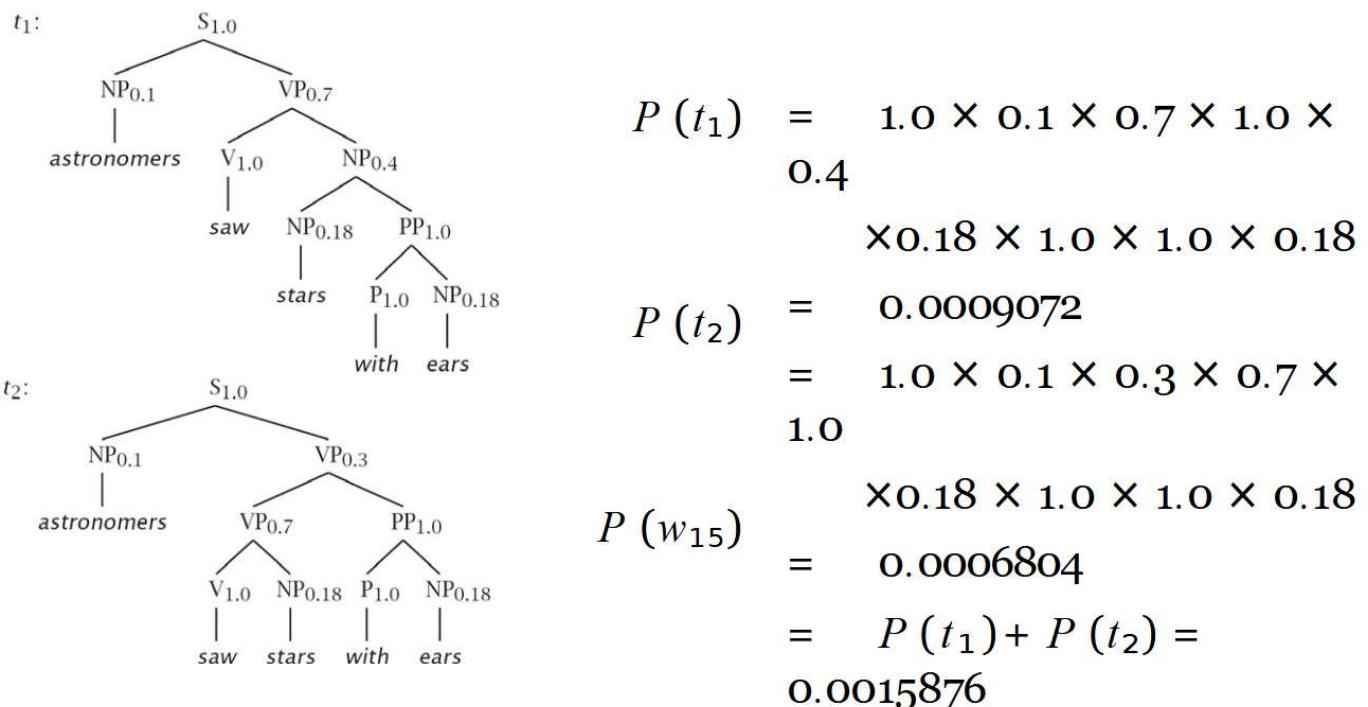
A PCFG consists of:

- A set of terminals, $\{w^k\}$, $k = 1, \dots, V$
- A set of nonterminals, N^i , $i = 1, \dots, n$
- A designated start symbol N^1
- A set of rules, $\{N^i \rightarrow \xi^j\}$, (where ξ^j is a sequence of terminals and nonterminals)
- A corresponding set of probabilities on rules such that: $\forall i \quad \sum_j P(N^i \rightarrow \xi^j) = 1$

$S \rightarrow NP VP$	1.0	$NP \rightarrow NP PP$	0.4
$PP \rightarrow P NP$	1.0	$NP \rightarrow astronomers$	0.1
$VP \rightarrow V NP$	0.7	$NP \rightarrow ears$	0.18
$VP \rightarrow VP PP$	0.3	$NP \rightarrow saw$	0.04
$P \rightarrow with$	1.0	$NP \rightarrow stars$	0.18
$V \rightarrow saw$	1.0	$NP \rightarrow telescopes$	0.1

Terminals with, saw, astronomers, ears, stars, telescopes

Nonterminals S, PP, P, NP, VP, V
 Start symbol S



Probabilistic CYK

	1	2	3	4	5
0	Det (.4)				
1					
2					
3					
4					
5					

- $S \rightarrow NP VP [.80]$
- $NP \rightarrow Det N [.30]$
- $VP \rightarrow V NP [.20]$
- $V \rightarrow includes [.05]$
- $Det \rightarrow the [.4]$
- $Det \rightarrow a [.4]$
- $N \rightarrow meal [.01]$
- $N \rightarrow flight [.02]$

The flight includes a meal.

	1	2	3	4	5
0	Det (.4)				
1		N .02			
2					
3					
4					
5					

- $S \rightarrow NP VP [.80]$
- $NP \rightarrow Det N [.30]$
- $VP \rightarrow V NP [.20]$
- $V \rightarrow includes [.05]$
- $Det \rightarrow the [.4]$
- $Det \rightarrow a [.4]$
- $N \rightarrow meal [.01]$
- $N \rightarrow flight [.02]$

The flight includes a meal.

	1	2	3	4	5
0	Det (.4)	NP .0024			
1		N .02			
2					
3					
4					
5					

- $S \rightarrow NP VP [.80]$
- $NP \rightarrow Det N [.30]$
- $VP \rightarrow V NP [.20]$
- $V \rightarrow includes [.05]$
- $Det \rightarrow the [.4]$
- $Det \rightarrow a [.4]$
- $N \rightarrow meal [.01]$
- $N \rightarrow flight [.02]$

The flight includes a meal.

Note: probability of NP in [0,2]
 $P(Det \rightarrow the) * P(N \rightarrow flight) * P(NP \rightarrow Det N)$

$$0.4 \times 0.02 \times 0.3 = 0.0024$$

	1	2	3	4	5
0	Det (.4)	NP .0024			
1		N .02			
2			V .05		
3					
4					
5					

- $S \rightarrow NP VP [.80]$
- $NP \rightarrow Det N [.30]$
- $VP \rightarrow V NP [.20]$
- $V \rightarrow includes [.05]$
- $Det \rightarrow the [.4]$
- $Det \rightarrow a [.4]$
- $N \rightarrow meal [.01]$
- $N \rightarrow flight [.02]$

The flight includes a meal.

$S \rightarrow NP VP [.80]$

$NP \rightarrow Det N [.30]$

$VP \rightarrow V NP [.20]$

$V \rightarrow includes [.05]$

$Det \rightarrow the [.4]$

$Det \rightarrow a [.4]$

$N \rightarrow meal [.01]$

$N \rightarrow flight [.02]$

$S \rightarrow NP VP [.80]$

$NP \rightarrow Det N [.30]$

$VP \rightarrow V NP [.20]$

$V \rightarrow includes [.05]$

$Det \rightarrow the [.4]$

$Det \rightarrow a [.4]$

$N \rightarrow meal [.01]$

$N \rightarrow flight [.02]$

	1	2	3	4	5
0	Det (.4)	NP .0024			
1		N .02			
2			V .05		
3				Det .4	
4					N .01
5					

The flight includes a meal.

	1	2	3	4	5
0	Det (.4)	NP .0024			
1		N .02			
2			V .05		
3				Det .4	N .01
4					
5					

The flight includes a meal.

	1	2	3	4	5
0	Det (.4)	NP .0024	X		
1		N .02	-X	X	
2			V .05	-X	Vf .0001
3			no production	Det .4	NP .001
4				{V Det}	N .01
5					no production

The flight includes a meal.

$P(Det \rightarrow a) \times P(N \rightarrow meal) \times P(NP \rightarrow Det N)$

$$0.4 \times 0.01 \times 0.3 = 0.0012$$

For cell [0,3] check pairs of non terms from cells [0,1] & [1,3] as [1,3] is empty check next
[0,2] & [2,3] is NP V there is no production

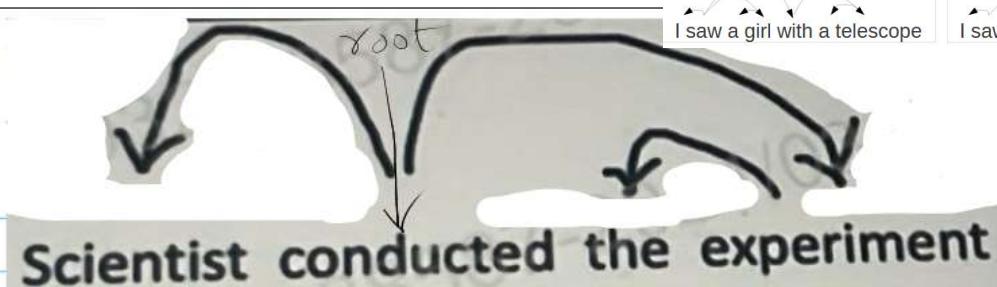
For cell [1,4] check pairs of non terms from cells [1,2] & [2,4] as [2,4] is empty check next
[1,3] & [3,4] as [1,3] is empty check next

For cell [2,5] check pairs of non terms from cells [2,3] & [3,5] is V NP exists VP $\rightarrow V NP$
 $P([2,5]) = P(V \rightarrow includes) \times P(N \rightarrow meal) \times P(VP \rightarrow V NP) = 0.5 \times 0.0012 \times 0.2 = 0.0001$

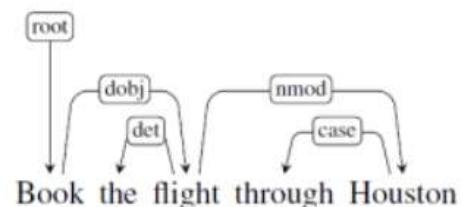
[2,4] & [4,5] as [2,4] is empty check next

If there were multiple probabilities then take $P([2,5]) = \max(P([2,5])_1, P([2,5])_2)$

Dependency Parsing/Arc Standard



transition	stack	Buffer	Action	Relation
0	[root]	[scientist, conducted, the, experiment]	shift	
1	[root, scientist]	[conducted, the, experiment]	shift	
2	[root, scientist, conducted]	[the, experiment]	left-Arc	scientist → conducted
3	[root, conducted]	[the, experiment]	shift	
4.	[root, conducted, the]	[experiment]	shift	
5.	[root, conducted, the, experiment]	[]	left-Arc	the → experiment
6.	[root, conducted, experiment]	[]	right-Arc	conduct → experiment
7	[root, conducted]	[]	right-Arc	
8	[root]	[]	Done	



Step	Stack	Word List	Action	Relation Added
0	[root]	[book, the, flight, through, houston]	RIGHTARC	(root → book)
1	[root, book]	[the, flight, through, houston]	SHIFT	
2	[root, book, the]	[flight, through, houston]	LEFTARC	(the ← flight)
3	[root, book]	[flight, through, houston]	RIGHTARC	(book → flight)
4	[root, book, flight]	[through, houston]	SHIFT	
5	[root, book, flight, through]	[houston]	LEFTARC	(through ← houston)
6	[root, book, flight]	[houston]	RIGHTARC	(flight → houston)
7	[root, book, flight, houston]	[]	REDUCE	
8	[root, book, flight]	[]	REDUCE	
9	[root, book]	[]	REDUCE	
10	[root]	[]	Done	

```

LEARN({ $T_1, \dots, T_N$ })
1    $w \leftarrow 0.0$ 
2   for  $i$  in  $1..K$ 
3     for  $j$  in  $1..N$ 
4        $c \leftarrow ([\emptyset_S, [w_1, \dots, w_{n_j}]_B, \{\}])$ 
5       while  $B_c \neq []$ 
6          $t^* \leftarrow \arg \max_t w \cdot f(c, t)$ 
7          $t_o \leftarrow o(c, T_i)$ 
8         if  $t^* \neq t_o$ 
9            $w \leftarrow w + f(c, t_o) - f(c, t^*)$ 
10           $c \leftarrow t_o(c)$ 
11   return  $w$ 

```

Consider the sentence, 'John saw Mary'.

Oracle $o(c, T_i)$ returns the optimal transition of c and T_i

Draw a dependency graph for this sentence.

Assume that you are learning a classifier for the data-driven deterministic parsing and the above sentence is a gold-standard parse in your training data. You are also given that *John* and *Mary* are 'Nouns', while the POS tag of *saw* is 'Verb'. Assume that your features correspond to the following conditions:

- C₀ ▪ The stack is empty
- C₁ ▪ Top of stack is Noun and Top of buffer is Verb
- C₂ ▪ Top of stack is Verb and Top of buffer is Noun



Initialize the weights of all your features to 5.0, except that in all of the above cases, you give a weight of 5.5 to *Left-Arc*. Define your feature vector and the initial weight vector.

Use this gold standard parse during online learning and report the weights after completing one full iteration of Arc-Eager parsing over this sentence.

Operations: Left-Arc LA, Right-Arc RA, Shift SH, Reduce RE

$$F(c,t) = [(c0,LA), (c1,LA), (c2,LA)] | [(c0,RA), (c1,RA), (c2,RA)] | [(c0,SH), (c1,SH), (c2,SH)] | [(c0,RE), (c1,RE), (c2,RE)]$$

$$W = [5.5, 5.5, 5.5, 5.5 | 5.0, 5.0, 5.0, 5.0 | 5.0, 5.0, 5.0, 5.0 | 5.0, 5.0, 5.0]$$

First iteration:

Stack	Buffer
[]	[John, Saw, Marry]

stack is empty which corresponds to feature c0 being true, thus

LA RA SH RE

$$F(c, LA) = [1, 0, 0 | 0, 0, 0 | 0, 0, 0 | 0, 0, 0]$$

$$F(c, RA) = [0, 0, 0 | 1, 0, 0 | 0, 0, 0 | 0, 0, 0]$$

$$F(c, SH) = [0, 0, 0 | 0, 0, 0 | 1, 0, 0 | 0, 0, 0]$$

$$F(c, RE) = [0, 0, 0 | 0, 0, 0 | 0, 0, 0 | 1, 0, 0]$$

$$t^* = \arg \max(W * F(c,t))$$

$$W * F(c, LA) = \text{sum}([5.5 * 1, 5.5 * 0, 5.5 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0]) \\ = \text{sum}([5.5, 0, 0, 0, 0, 0, 0, 0, 0, 0]) = 5.5$$

$$W * F(c, RA) = \text{sum}([5.5 * 0, 5.5 * 0, 5.5 * 0 | 5.0 * 1, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0]) \\ = \text{sum}([0, 0, 0, 5.0, 0, 0, 0, 0, 0, 0]) = 5.0$$

$$W * F(c, SH) = \text{sum}([5.5 * 0, 5.5 * 0, 5.5 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0]) \\ = \text{sum}([0, 0, 0, 0, 0, 5.0, 0, 0, 0, 0]) = 5.0$$

$$W * F(c, RE) = \text{sum}([5.5 * 0, 5.5 * 0, 5.5 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 0, 5.0 * 0, 5.0 * 0 | 5.0 * 1, 5.0 * 0, 5.0 * 0]) \\ = \text{sum}([0, 0, 0, 0, 0, 0, 5.0, 0, 0]) = 5.0$$

It was LA operation which gave us the maximum value of 5.5 therefore $t^* = LA$

As per oracle the optimal transition t_0 is SH (when $\text{len(stack)} > 0$ else if $\text{len(stack)} == 0$ then t_0 is RE)

Updating weights

$$W = W + F(c, t_0) - F(c, t^*)$$

$$= [5.5, 5.5, 5.5, 5.5 | 5.0, 5.0, 5.0, 5.0 | 5.0, 5.0, 5.0, 5.0] + [0, 0, 0 | 0, 0, 0 | 1, 0, 0 | 0, 0, 0] - [1, 0, 0 | 0, 0, 0 | 0, 0, 0 | 0, 0, 0]$$

$$W = [4.5, 5.5, 5.5, 5.5 | 5.0, 5.0, 5.0, 5.0 | 6.0, 5.0, 5.0, 5.0 | 5.0, 5.0, 5.0]$$

Iteration 1 done.

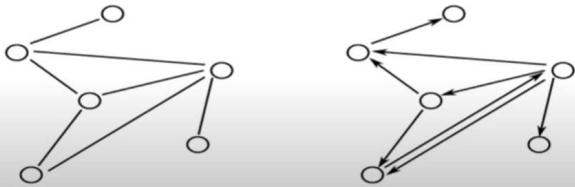
Graph-Based Dependency Parsing

Some Graph Theory Reminders

- A graph $G = (V, A)$ is a set of vertices V and arcs $(i, j) \in A$ where $i, j \in V$.
- Undirected graphs: $(i, j) \in A \Leftrightarrow (j, i) \in A$
- Directed graphs (digraphs): $(i, j) \in A \Rightarrow (j, i) \notin A$

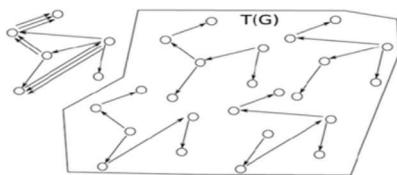
Multi-Digraphs

- A multi-digraph is a digraph where multiple arcs between vertices are possible
- $(i, j, k) \in A$ represents the k^{th} arc from vertex i to vertex j .



Maximum Spanning Trees (MST)

Let $T(G)$ be the set of all spanning trees for graph G

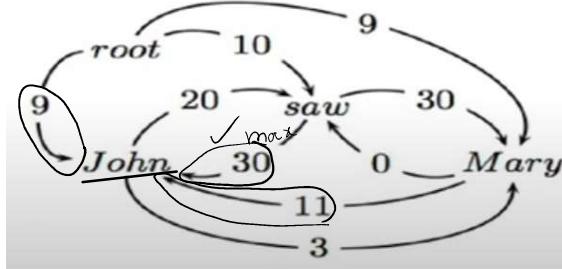
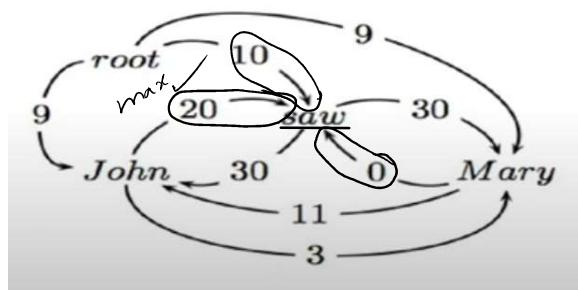
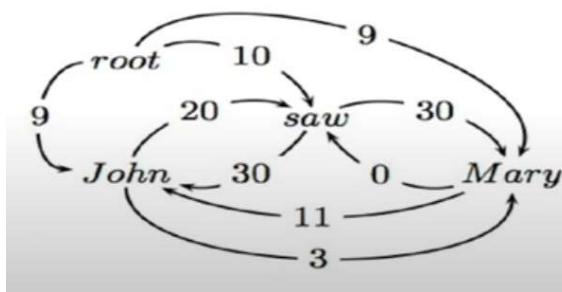


The MST problem

Find the spanning tree G' of the graph G that has the highest weight

$$G' = \arg \max_{G' \in T(G)} w(G') = \arg \max_{G' \in T(G)} \sum_{(i,j,k) \in G'} w_{ij}^k$$

$x = \text{John saw Mary}$



Chu-Liu-Edmonds algorithm

function MAXSPANNINGTREE($G=(V,E)$, root , score) returns spanning tree

```

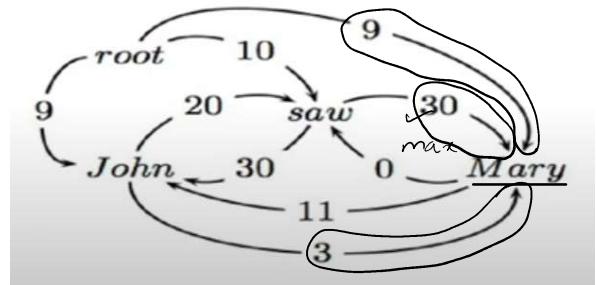
 $F \leftarrow []$ 
 $T' \leftarrow []$ 
 $\text{score}' \leftarrow []$ 
for each  $v \in V$  do
     $\text{bestInEdge} \leftarrow \arg\max_{e=(u,v) \in E} \text{score}[e]$ 
     $F \leftarrow F \cup \text{bestInEdge}$ 
    for each  $e=(u,v) \in E$  do
         $\text{score}'[e] \leftarrow \text{score}[e] - \text{score}[\text{bestInEdge}]$ 
if  $T=(V,F)$  is a spanning tree then return it
else
     $C \leftarrow$  a cycle in  $F$ 
     $G' \leftarrow \text{CONTRACT}(G, C)$ 
     $T' \leftarrow \text{MAXSPANNINGTREE}(G', \text{root}, \text{score}')$ 
     $T \leftarrow \text{EXPAND}(T', C)$ 
return  $T$ 

```

[select best incoming edge for each node]
 [subtract its score from all incoming edges]
 [stopping condition]
 [contract nodes if there are cycles]
 [recursively compute MST]
 [expand contracted nodes]

function CONTRACT(G, C) returns contracted graph

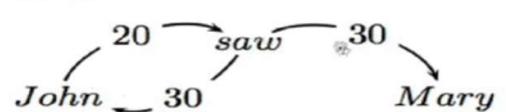
function EXPAND(T, C) returns expanded graph



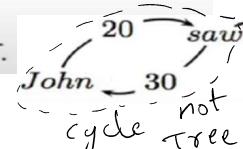
Find the highest scoring incoming arc for each vertex

no incoming edges to root

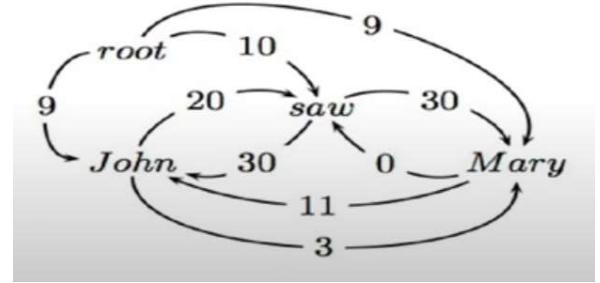
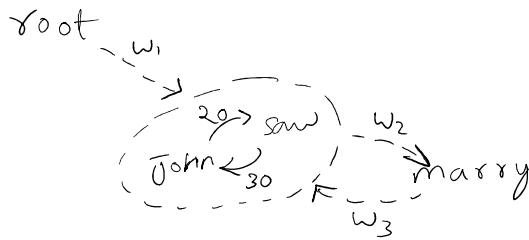
root



If this is a tree, then we have found MST.

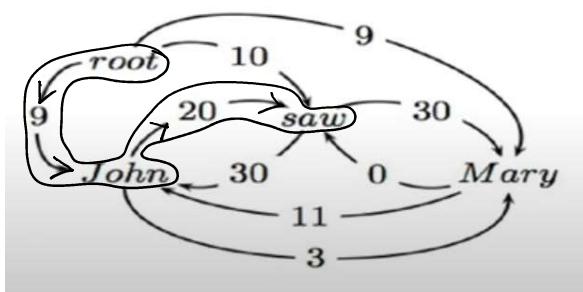


since there is cycle contract
i.e make nodes John & saw
as one new temporary node



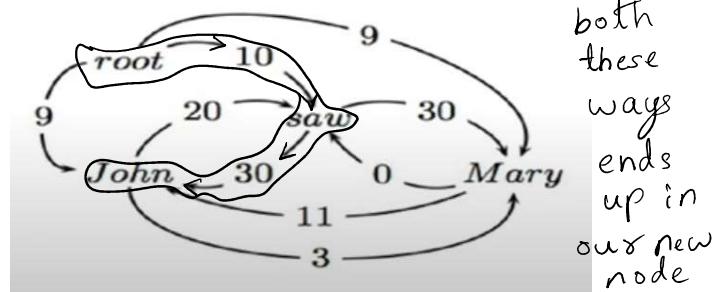
since it's a new node it may have incoming & outgoing edges to every other non-root nodes.

if we consider unique paths b/w root, John, saw
with no cycles & in single direction there are 2 ways



$$\text{root} \xrightarrow{9} \text{John} \xrightarrow{20} \text{saw}$$

$\underbrace{\hspace{2cm}}_{29}$



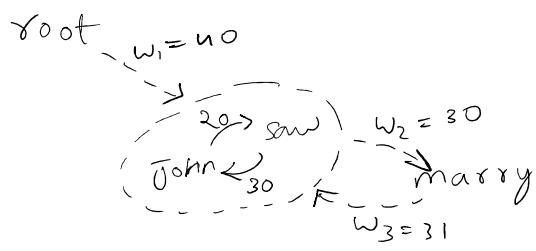
$$\text{root} \xrightarrow{10} \text{saw} \xrightarrow{30} \text{John}$$

$\underbrace{\hspace{2cm}}_{40}$

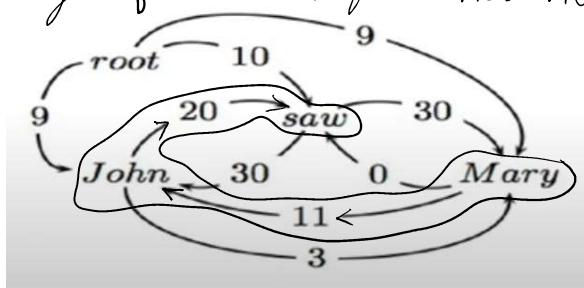
both these ways ends up in our new node

$$\text{weight for root to new node } w_1 = \max(29, 40) = 40$$

to find w_2 i.e outgoing edge from irrespective of unique paths there are only 2 exit points from this

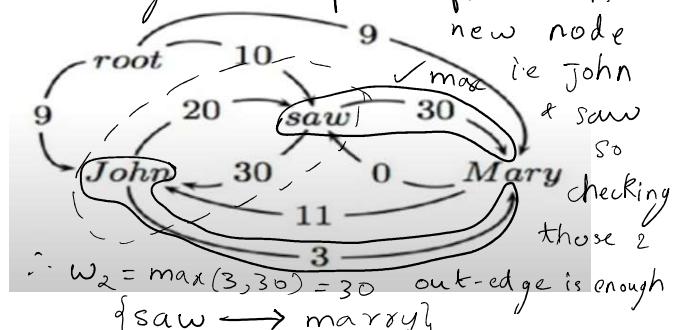


to find w_3 weight of incoming edge from marry to new node

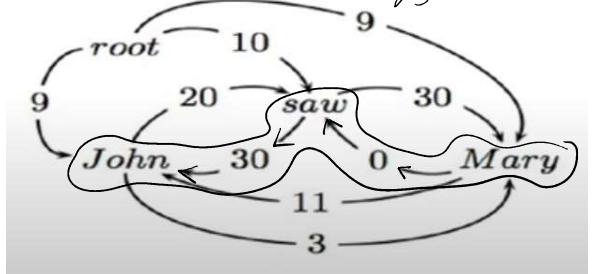


$$\text{marry} \xrightarrow{11} \text{John} \xrightarrow{20} \text{saw}$$

$$w_3 = \max(30, 31) = 31$$



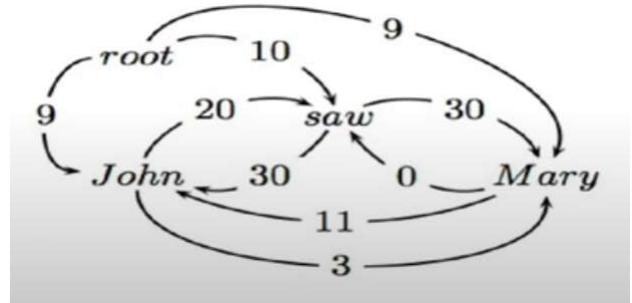
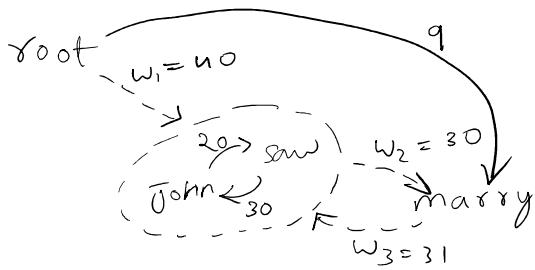
$$\therefore w_2 = \max(30, 30) = 30 \text{ out-edge is enough } \{ \text{saw} \rightarrow \text{marry} \}$$



$$\text{marry} \xrightarrow{0} \text{saw} \xrightarrow{30} \text{John}$$

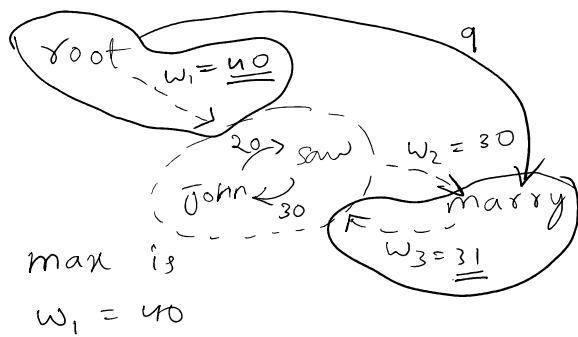
$\underbrace{\hspace{2cm}}_{30}$

It's still has cycles
thus not a tree yet

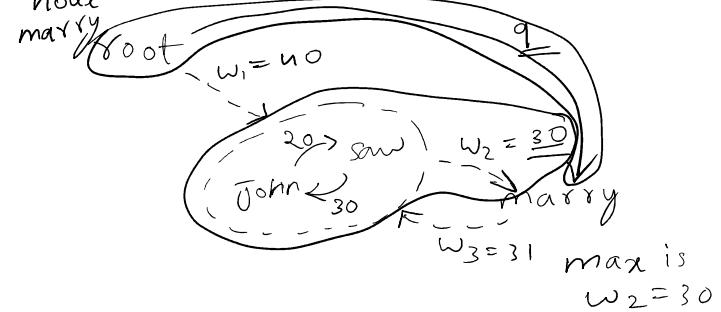


To get max spanning tree MST choose max weighted incoming edge for every node

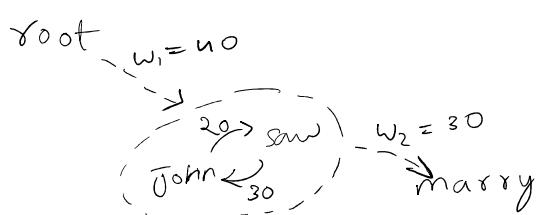
for
new
node



for
node
mary



after preserving max edges



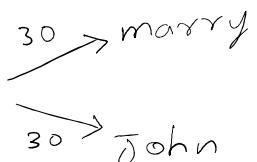
In
no further contraction possible

we still have 1 contracted node
so start expanding

by combining original edges we get
our final dependency graph

$w_2 = 30$ is contributed by
original edge $\text{saw} \xrightarrow{30} \text{marry}$

$w_1 = u0$ contributed by
 $\text{root} \xrightarrow{10} \text{saw} \xrightarrow{30} \text{John}$



Graph-Based MST Learning

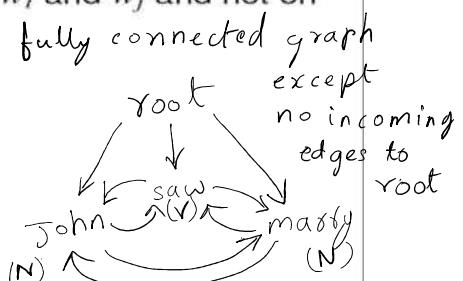
Training data: $T = \{(x_t, G_t)\}_{t=1}^{|T|}$ G_t : Gold standard parse tree

1. $w^{(0)} = 0; i = 0$
2. for $n : 1..N$
3. for $t : 1..|T|$
4. Let $G' = \text{argmax}_{G'} w^{(i)}.f(G')$
5. if $G' \neq G_t$
6. $w^{(i+1)} = w^{(i)} + f(G_t) - f(G')$
7. $i = i + 1$
8. return w^i

Suppose you are training MST Parser for dependency and the sentence, "John saw Mary" occurs in the training set. Also, for simplicity, assume that there is only one dependency relation, "rel". Thus, for every arc from word w_i to w_j , your features may be simplified to depend only on words w_i and w_j and not on the relation label.

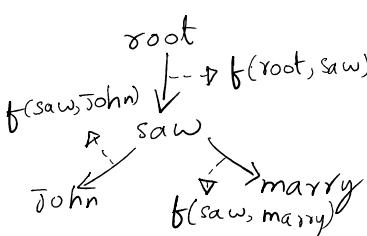
Below is the set of features

- $f_1: \text{pos}(w_i) = \text{Noun}$ and $\text{pos}(w_j) = \text{Noun}$
- $f_2: \text{pos}(w_i) = \text{Verb}$ and $\text{pos}(w_j) = \text{Noun}$
- $f_3: w_i = \text{Root}$ and $\text{pos}(w_j) = \text{Verb}$
- $f_4: w_i = \text{Root}$ and $\text{pos}(w_j) = \text{Noun}$
- $f_5: w_i = \text{Root}$ and w_j occurs at the end of sentence
- $f_6: w_i$ occurs before w_j in the sentence
- $f_7: \text{pos}(w_i) = \text{Noun}$ and $\text{pos}(w_j) = \text{Verb}$



The feature weights before the start of the iteration are: {3, 20, 15, 12, 1, 10, 20}. Determine the weights after an iteration over this example.

let G_t Gold standard parse tree



from chu-Liu edmond
we got G'

Edge $f(w_i, w_j)$	f_1	f_2	f_3	f_4	f_5	f_6	f_7
$f(\text{root}, \text{john})$	0	0	0	1	0	0	0
$f(\text{root}, \text{saw})$	0	0	1	0	0	0	0
$f(\text{root}, \text{marry})$	0	0	0	1	1	0	0
$f(\text{john}, \text{saw})$	0	0	0	0	0	1	1
$f(\text{john}, \text{marry})$	1	0	0	0	0	1	0
$f(\text{saw}, \text{john})$	0	1	0	0	0	0	0
$f(\text{saw}, \text{marry})$	0	1	0	0	0	1	0
$f(\text{marry}, \text{john})$	1	0	0	0	0	0	0
$f(\text{marry}, \text{saw})$	0	0	0	0	0	0	1

$$w = [3, 20, 15, 12, 1, 10, 20]$$

Edge weight = dot product of w & $f(w_i, w_j)$

for edge ('root', 'john') weight: 12

for edge ('root', 'saw') weight: 15

for edge ('root', 'marry') weight: 13

for edge ('john', 'saw') weight: 30

for edge ('john', 'marry') weight: 13

for edge ('saw', 'john') weight: 20

for edge ('saw', 'marry') weight: 30

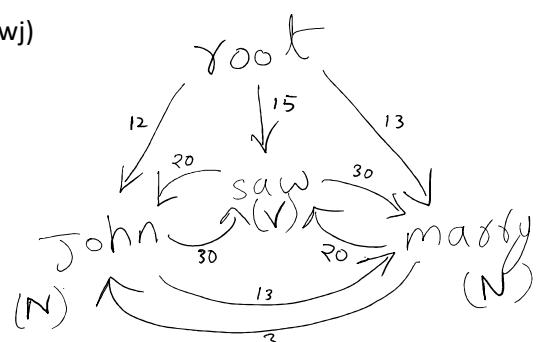
for edge ('marry', 'john') weight: 3

for edge ('marry', 'saw') weight: 20

$$f(G_t) = f(\text{root}, \text{saw}) + f(\text{saw}, \text{john}) + f(\text{saw}, \text{marry})$$

$$f(\text{saw}, \text{marry}) = [0, 2, 1, 0, 0, 1, 0]$$

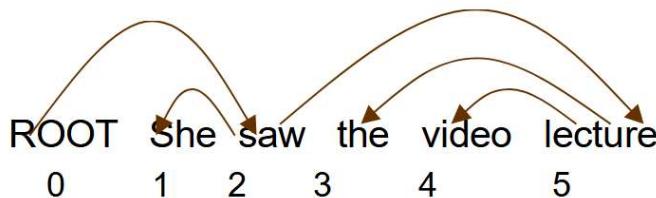
$$f(G') = f(\text{'root'}, \text{'john'}) + f(\text{'john'}, \text{'saw'}) + f(\text{'john'}, \text{'marry'}) = [1, 0, 0, 1, 0, 2, 1]$$



$$w_{\text{new}} = w + f(G_t) - f(G')$$

$$= [2, 22, 16, 11, 1, 9, 19]$$

- Unlabeled Attachment Score (UAS), which corresponds to the number of correctly predicted dependencies over the number of possibilities;
 - Labeled Attachment Score (LAS), which corresponds to the number of correctly predicted dependencies and relations over the number of possibilities.



$$\text{Acc} = \frac{\# \text{ correct deps}}{\# \text{ of deps}}$$

$$\text{UAS}(\text{unlabeled attachment score}) = 4 / 5 \\ = 80\%$$

$$\text{LAS} (\text{ Labeled Attachment Score}) = 2 / 5 \\ = 40\%$$

Gold			
1	2	She	nsubj
2	0	saw	root
3	5	the	det
4	5	video	nn
5	2	lecture	dobj

Parsed			
1	2	She	nsubj
2	0	saw	root
3	4	the	det
4	5	video	nsubj
5	2	lecture	ccomp

UAS
LAS
out of no. of UAS

how many matching relations

Word Sense Disambiguation

What's a word sense?

- A sense (or word sense) is a discrete representation of one aspect of the meaning of a word
- Example
 - We went to see the play Romeo and Juliet at the theater.
 - The children went out to play in the park.

Homonymy

Lexemes that share a form

- Phonological, orthographic or both

But have unrelated, distinct meanings

- Examples
 - bat (wooden stick-like thing) vs bat (flying scary mammal thing)
 - bank (financial institution) versus bank (riverside)
- Can be homophones, homographs, or both:
 - Homophones:
 - Write and right
 - Piece and peace

Polysemy

The **bank** is constructed from red brick

I withdrew the money from the **bank**

- Which sense of bank is this?

- Is it distinct from the river bank sense?
- How about the savings bank sense?

Another example:

His cottage is near a small **wood**.

The statue was made out of a block of **wood**.

Are those the same sense?

Synonyms

Words that have the same meaning in some or all contexts.

- filbert / hazelnut
- couch / sofa
- big / large
- automobile / car
- vomit / throw up
- Water / H₂O

Two lexemes are synonyms if they can be successfully substituted for each other in all situations

Relationships between word senses

- Homonymy
- Polysemy
- Synonymy
- Antonymy
- Hyponymy
- Hypernymy

Homonymy causes problems for NLP applications

- Text-to-Speech
 - Bass vs bass
- Information retrieval
 - Bat care
- Machine Translation
 - I went to bank
- Speech recognition
 - Piece and peace

A single lexeme with multiple **related** meanings (bank the building, bank the financial institution)

Most non-rare words have multiple meanings

- The number of meanings is related to its frequency
- Verbs tend more to polysemy
- Distinguishing polysemy from homonymy isn't always easy (or necessary)

But

There are no examples of perfect synonymy

- Why should that be?
- Even if many aspects of meaning are identical
- Still may not preserve the acceptability based on notions of politeness, slang, register, genre, etc.

Example:

- Water and H₂O

Word sense disambiguation (contd.)

Synonymy is a relation between senses rather than words

Consider the words **big** and **large**

Are they synonyms?

- How **big** is that plane?
- Would I be flying on a **large** or small plane?

How about here:

- Miss Nelson, for instance, became a kind of **big** sister to Benjamin.
- ?Miss Nelson, for instance, became a kind of **large** sister to Benjamin.

Why?

- **big** has a sense that means being older, or grown up
- **large** lacks this sense

Hyponymy

- One sense is a **hyponym** of another if the first sense is more specific, denoting a subclass of the other
 - *car* is a hyponym of *vehicle*
 - *dog* is a hyponym of *animal*
 - *mango* is a hyponym of *fruit*
- Conversely
 - *vehicle* is a hypernym/superordinate of *car*
 - *animal* is a hypernym of *dog*
 - *fruit* is a hypernym of *mango*

superordinate	vehicle	fruit	furniture	mammal
hyponym	car	mango	chair	dog

WordNet Noun Relations

Relation	Also called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> ¹ → <i>meal</i> ¹
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> ¹ → <i>lunch</i> ¹
Member Meronym	Has-Member	From groups to their members	<i>faculty</i> ² → <i>professor</i> ¹
Has-Instance		From concepts to instances of the concept	<i>composer</i> ¹ → <i>Bach</i> ¹
Instance		From instances to their concepts	<i>Austen</i> ¹ → <i>author</i> ¹
Member Holonym	Member-Of	From members to their groups	<i>copilot</i> ¹ → <i>crew</i> ¹
Part Meronym	Has-Part	From wholes to parts	<i>table</i> ² → <i>leg</i> ³
Part Holonym	Part-Of	From parts to wholes	<i>course</i> ⁷ → <i>meal</i> ¹
Antonym		Opposites	<i>leader</i> ¹ → <i>follower</i> ¹

WordNet Verb Relations

Relation	Definition	Example
Hypernym	From events to superordinate events	<i>fly</i> ⁹ → <i>travel</i> ³
Troponym	From a verb (event) to a specific manner elaboration of that verb	<i>walk</i> ¹ → <i>stroll</i> ¹
Entails	From verbs (events) to the verbs (events) they entail	<i>snore</i> ¹ → <i>sleep</i> ¹
Antonym	Opposites	<i>increase</i> ¹ ↔ <i>decrease</i> ¹

Antonyms

Senses that are opposites with respect to one feature of their meaning

Otherwise, they are very similar!

- dark / light
- short / long
- hot / cold
- up / down
- in / out

More formally: antonyms can

- define a binary opposition or at opposite ends of a scale (*long/short*, *fast/slow*)
- Be **reverses**: *rise/fall*, *up/down*

Wordnet

- The set of near-synonyms for a WordNet sense is called a **synset (synonym set)**; it's their version of a sense or a concept

Example: **chump** as a noun to mean

'a person who is gullible and easy to take advantage of'

{*chump*¹, *fool*², *gull*¹, *mark*⁹, *patsy*¹, *fall guy*¹, *sucker*¹,
*soft touch*¹, *mug*²}

- Each of these senses share this same gloss
- Thus for WordNet, the meaning of this sense of **chump** is this list.

WordNet Hierarchies

Sense 3
bass, *basso* --
 (an adult male singer with the lowest voice)
 => *singer*, *vocalist*, *vocalizer*, *vocaliser*
 => *musician*, *instrumentalist*, *player*
 => *performer*, *performing artist*
 => *entertainer*
 => *person*, *individual*, *someone...*
 => *organism*, *being*
 => *living thing*, *animate thing*,
 => *whole*, *unit*
 => *object*, *physical object*
 => *physical entity*
 => *entity*
 => *causal agent*, *cause*, *causal agency*
 => *physical entity*
 => *entity*

Sense 7
bass--
 (the member with the lowest range of a family of musical instruments)
 => *musical instrument*, *instrument*
 => *device*
 => *instrumentality*, *instrumentation*
 => *artifact*, *artefact*
 => *whole*, *unit*
 => *object*, *physical object*
 => *physical entity*
 => *entity*

Word sense disambiguation (contd.)

Word Sense Disambiguation (WSD)

The task of selecting the correct sense for a word is called word sense disambiguation, or WSD

Given

- a word in context,
- a fixed inventory of potential word sense

Decide which sense of the word this is

Examples

- English-to-Spanish MT
 - Inventory is set of Spanish translations
- Speech Synthesis
 - Inventory is homographs with different pronunciations like *bass* and *bow*

Two variants of WSD task

Lexical sample task

- Small pre-selected set of target words
- And inventory of senses for each word
- Since these words and the set of senses are small, simple supervised classification approaches work very well

All-words task

- In this all-words task, the system is given an all-words entire texts and
- lexicon with an inventory of senses for each entry
- we have to disambiguate every word in the text (or sometimes just every content word).

WSD Methods Supervised Machine Learning Approaches

- Supervised
- Dictionary based
- Semi supervised

Supervised machine learning approach:

- a **training corpus** of words tagged in context with their sense
- used to train a classifier that can tag words in new text

Summary of what we need:

- the **tag set** ("sense inventory")
- the **training corpus**
- A set of **features** extracted from the training corpus
- A **classifier**

Supervised WSD 1: WSD Tags

What's a tag?

- A dictionary sense?

For example, for WordNet an instance of "bass" in a text has 8 possible tags or labels (bass1 through bass8).

WordNet Bass

The noun ``bass'' has 8 senses in WordNet

1. bass - (the lowest part of the musical range)
2. bass, bass part - (the lowest part in polyphonic music)
3. bass, basso - (an adult male singer with the lowest voice)
4. sea bass, bass - (flesh of lean-fleshed saltwater fish of the family Serranidae)
5. freshwater bass, bass - (any of various North American lean-fleshed freshwater fishes especially of the genus Micropterus)
6. bass, bass voice, basso - (the lowest adult male singing voice)
7. bass - (the member with the lowest range of a family of musical instruments)
8. bass - (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

Word sense disambiguation (contd.)

Supervised WSD 2: Get a corpus

innovate achieve

Lexical sample task:

- Line-hard-serve corpus - 4000 examples of each
- Interest corpus - 2369 sense-tagged examples

All words:

- **Semantic concordance**: a corpus in which each open-class word is labeled with a sense from a specific dictionary/thesaurus.
 - SemCor: 234,000 words from Brown Corpus, manually tagged with WordNet senses
 - SENSEVAL-3 competition corpora - 2081 tagged word tokens

Example of SemCor with wordnet sense numbers

You will find⁹ *v* that avocado¹ *n* is¹ *v* unlike¹ *j* other¹ *j* fruit¹ *n* you have ever¹ *r*, tasted² *v*

Given each noun, verb, adjective, or adverb word in the hand-labeled test set. Ex: fruit¹ *n* (the ripened reproductive body of a seed plant), and the other two senses fruit² *n* (yield;an amount of a product) and fruit³ *n* (the consequence of some effort or action).

Supervised WSD 3: Extract feature vectors

A simple representation for each observation (each instance of a target word)

- Vectors of sets of feature/value pairs
 - i.e. files of comma-separated values
- These vectors should represent the window of words around the target

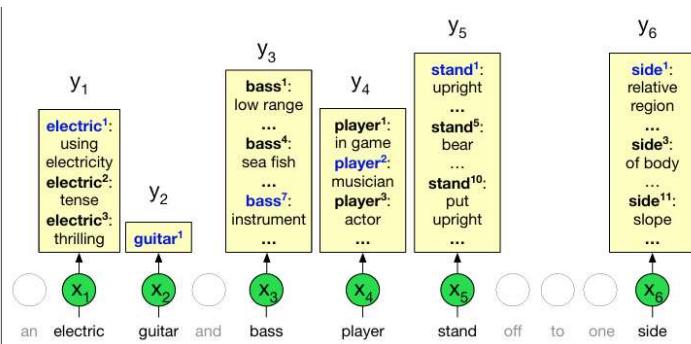


Figure 23.8 The all-words WSD task, mapping from input words (x) to WordNet senses (y). Only nouns, verbs, adjectives, and adverbs are mapped, and note that some words (like guitar in the example) only have one sense in WordNet. Figure inspired by Chaplot and Salakhutdinov (2018).

Two kinds of features in the vectors

Collocational features and bag-of-words features

– Collocational

- Features about words at **specific** positions near target word
 - Often limited to just word identity and POS

– Bag-of-words

- Features about words that occur anywhere in the window (regardless of position)
 - Typically limited to frequency counts

Collocational

Position-specific information about the words in the window

guitar and bass player stand

- [guitar, NN, and, CC, player, NN, stand, VB]
- Word_{n-2}, POS_{n-2}, word_{n-1}, POS_{n-1}, Word_{n+1} POS_{n+1}...
- In other words, a vector consisting of
- [position n word, position n part-of-speech...]

Example text (WSJ)

An electric **guitar** and **bass** **player** **stand** off to one side not really part of the scene,

Assume a window of +/- 2 from the target

Word sense disambiguation (contd.)

Bag-of-words

Words that occur within the window, regardless of specific position

First derive a set of terms to place in the vector

Then note how often each of those terms occurs in a given window

Assume we've settled on a possible vocabulary of 12 words that includes **guitar** and **player** but not **and** and **stand**

[*fishing, big, sound, player, fly, rod, pound, double, runs, playing, guitar, band*]

- The vector for:

guitar **and** **bass** **player** **stand**

[0,0,0,1,0,0,0,0,0,0,1,0]

Supervised WSD4:Classifier

Any kind of classifier

Input:

- a word **w** in a text window **d** (which we'll call a "document")
- a fixed set of classes $C = \{c_1, c_2, \dots, c_n\}$
- A training set of **m** hand-labeled text windows again called "documents" $(d_1, c_1), \dots, (d_m, c_m)$

- Naive Bayes
- Logistic regression
- Neural Networks
- Support-vector machines
- k-Nearest Neighbors

Output:

- a learned classifier $y:d \rightarrow c$

Lesk Algorithm Example

innovate achieve lead

- Consider disambiguating the word **bank** in the following context:

The bank can guarantee deposits will eventually cover future tuition costs because it invests in adjustable-rate mortgage securities.

- Wordnet senses of "bank":

bank ¹	Gloss:	a financial institution that accepts <u>deposits</u> and channels the money into lending activities
	Examples:	"he cashed a check at the bank", "that bank holds the <u>mortgage</u> on my home"
bank ²	Gloss:	sloping land (especially the slope beside a body of water)
	Examples:	"they pulled the canoe up on the bank", "he sat on the bank of the river and watched the currents"

- Sense bank1 has two non-stopwords overlapping with the context in deposits and mortgage, while sense bank2 has zero words, so bank1 is chosen.

Word Sense Disambiguation Naïve Bayes Example

Using the Naïve Bayes classifier, we need to disambiguate the word "bank" in the Test Document #1 given in the table below. The Training Documents and their respective classifications are given below. [6 Marks]

	Doc	Sentence	Tokenized Words	Class
Test	1	I went to the bank to deposit my money.	Bank, deposit, money	?
Training	2	The fish swam in the river bank.	Fish, swim, river, bank	R=River_Bank
Training	3	We need to take out a loan from the bank	Loan, bank	F=Financial_Institution
Training	4	The boat was moored at the river bank	Boat, river, bank	R=River_Bank
Training	5	I deposited some money in the bank	deposit, money, bank	F=Financial_Institution

$$P(R) = \frac{2}{4} = 0.5 \quad \begin{matrix} \text{no. of unique words in} \\ \text{vocabulary} = 8 \end{matrix}$$

$$P(F) = 2/4 = 0.5 \quad \begin{matrix} \text{no. of words in} \\ \text{class } R = 5 \\ \text{class } F = 5 \end{matrix}$$

Test: Bank, deposit, money

$$P(\text{Bank} | R) = \frac{2+1}{5+8} = \frac{3}{13} = 0.23$$

$$P(\text{deposit} | R) = \frac{0+1}{5+8} = \frac{1}{13} = 0.0769$$

$$P(\text{money} | R) = \frac{0+1}{5+8} = \frac{1}{13} = 0.0769$$

$$P(\text{Bank} | F) = \frac{2+1}{5+8} = \frac{3}{13} = 0.23$$

$$P(\text{deposit} | F) = \frac{1+1}{5+8} = \frac{2}{13} = 0.1538$$

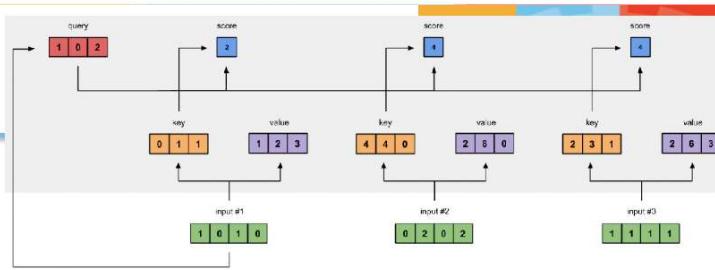
$$P(\text{money} | F) = \frac{1+1}{5+8} = \frac{2}{13} = 0.1538$$

$$\begin{aligned} P(\text{Test} | R) &\propto P(R) \times P(\text{Bank} | R) \times P(\text{deposit} | R) \times P(\text{money} | R) \\ &= 0.5 \times 0.23 \times (0.0769)^2 = 0.00068 \end{aligned}$$

$$\begin{aligned} P(\text{Test} | F) &\propto P(F) \times P(\text{Bank} | F) \times P(\text{deposit} | F) \times P(\text{money} | F) \\ &= 0.5 \times 0.23 \times (0.1538)^2 = 0.0027 \end{aligned}$$

as $P(\text{Test} | F) > P(\text{Test} | R)$ the test doc belongs to class financial_institution (F)

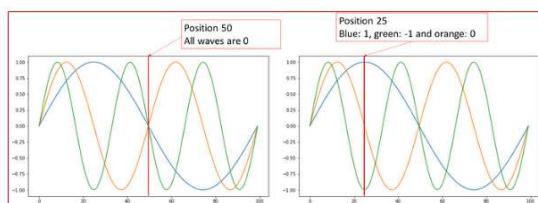
Self attention Example

		3 input vectors	
		$q_i = Qx_i$ (queries)	
		$k_i = Kx_i$ (keys)	
		$v_i = Vx_i$ (values)	
Weights for Key		Key representation for Input	Key for all inputs
$[[0, 0, 1],$ $[1, 1, 0],$ $[0, 1, 0],$ $[1, 1, 0]]$		$[0, 0, 1]$ $[1, 0, 1] \cdot [1, 1, 0] = [0, 1, 1]$ $[0, 1, 0]$ $[1, 1, 0]$	$[0 \ 1 \ 1]$ $[4 \ 4 \ 0]$ $[2 \ 3 \ 1]$
Weights for query		Query representation for Input	Query for all inputs
$[[1, 0, 1],$ $[1, 0, 0],$ $[0, 0, 1],$ $[0, 1, 1]]$		$[1, 0, 1]$ $[1, 0, 1] \cdot [1, 0, 0] = [1, 0, 2]$ $[0, 0, 1]$ $[0, 1, 1]$	$[1, 0, 2]$ $[2, 2, 2]$ $[2, 1, 3]$
Weights for Value		Value representation for Input	Value for all inputs
$[[0, 2, 0],$ $[0, 3, 0],$ $[1, 0, 3],$ $[1, 1, 0]]$		$[0, 2, 0]$ $[1, 0, 1] \cdot [0, 3, 0] = [1, 2, 3]$ $[1, 0, 3]$ $[1, 1, 0]$	$[1, 2, 3]$ $[2, 8, 0]$ $[2, 6, 3]$
CONTD...			
$e_{ij} = q_i^T k_j$ $\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$			
3 attention scores for first input i : (e_{ij})			
$[0, 1, 1]$			
$[1, 0, 2] \cdot [4, 4, 0] = [2, 4, 4]$			
$[2, 3, 1]$			
3 softmax attention scores (α_{ij})		Output 1, which is based on the query representation from Input 1 interacting with all other keys, including itself	
$\text{softmax}([2, 4, 4]) = [0.0, 0.5, 0.5]$		$[0.0, 0.0, 0.0]$ $+ [1.0, 4.0, 0.0]$ $+ [1.0, 3.0, 1.5]$ ----- $= [2.0, 7.0, 1.5]$	
weighted values			
1: $0.0 * [1, 2, 3] = [0.0, 0.0, 0.0]$		All the outputs after first iteration	
2: $0.5 * [2, 8, 0] = [1.0, 4.0, 0.0]$		$[2.0, 7.0, 1.5]$, # Output 1	
3: $0.5 * [2, 6, 3] = [1.0, 3.0, 1.5]$		$[2.0000, 8.0, 0.0]$, # Output 2	
		$[2.0, 7.8, 0.3]$ # Output 3	

Position representation vectors through sinusoids

Sinusoidal position representations: concatenate sinusoidal functions of varying periods

$$\mathbf{p}_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



Sequence	Index of token, k	Positional Encoding Matrix with $d=4$, $n=100$			
		$i=0$	$i=0$	$i=1$	$i=1$
I	0	$P_{00}=\sin(0) = 0$	$P_{01}=\cos(0) = 1$	$P_{02}=\sin(0) = 0$	$P_{03}=\cos(0) = 1$
am	1	$P_{10}=\sin(1/1) = 0.84$	$P_{11}=\cos(1/1) = 0.54$	$P_{12}=\sin(1/10) = 0.10$	$P_{13}=\cos(1/10) = 1.0$
a	2	$P_{20}=\sin(2/1) = 0.91$	$P_{21}=\cos(2/1) = -0.42$	$P_{22}=\sin(2/10) = 0.20$	$P_{23}=\cos(2/10) = 0.98$
Robot	3	$P_{30}=\sin(3/1) = 0.14$	$P_{31}=\cos(3/1) = -0.99$	$P_{32}=\sin(3/10) = 0.30$	$P_{33}=\cos(3/10) = 0.96$

Positional Encoding Matrix for the sequence 'I am a robot'

WordPiece tokenization algorithm

WordPiece is the tokenization algorithm Google developed to pretrain BERT.

Starts from a small vocabulary : special tokens used by the model and the initial alphabet.

Identifies subwords by adding a prefix (like ## for BERT),

"word" gets split like this: w ##o ##r ##d

Initial vocabulary contains all the characters present at the beginning of a word and the characters present inside a word preceded by the WordPiece prefix.

WordPiece tokenization algorithm

Instead of selecting the most frequent pair (BPE), WordPiece computes a score for each pair

$$\text{score} = (\text{freq_of_pair}) / (\text{freq_of_first_element} \times \text{freq_of_second_element})$$

Algorithm prioritizes the merging of pairs where the individual parts are less frequent in the vocabulary.

- It won't necessarily merge ("un", "#able") even if that pair occurs very frequently in the vocabulary, because the two pairs "un" and "#able" will likely each appear in a lot of other words and have a high frequency.
- ("hu", "#gging") will probably be merged faster (assuming the word "hugging" appears often in the vocabulary) since "hu" and "#gging" are likely to be less frequent individually.

Tokenization example using wordpiece

Corpus uses these five words: "hug", "pug", "pun", "bun", "hugs"

Assume the words had the following frequencies: ("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)

Training: token **learner** : splitting words and adding a prefix ## for subwords

- ("h" "#u" "#g", 10), ("p" "#u" "#g", 5), ("p" "#u" "#n", 12),
- ("b" "#u" "#n", 4), ("h" "#u" "#g" "#s", 5)

First iteration:

initial vocabulary : ["b", "h", "p", "#g", "#n", "#s", "#u"] (ignoring special tokens)

- Score ("#u", "#g") : pair freq/individual freq = $20/36*20=1/36$
- Score ("#g", "#s") : pair freq/individual freq = $5/20*5=1/20$
- Likewise compute scores for all the pairs
- first merge learned is ("#g", "#s") -> (**#gs**)

Vocabulary after first iteration: ["b", "h", "p", "#g", "#n", "#s", "#u", "**#gs**"]

Repeat iterations till you get a desired vocab size.

Token **segmenter**: New inputs are tokenized by applying the following steps:

- Pre-tokenize it
- Split it
- Apply the tokenization algorithm on each word.
 - Look for the biggest subword starting at the beginning of the first word and split it
 - Repeat the process on the second part
 - And so on for the rest of that word and the following words in the text

Continuing the example: suppose the final vocab learnt:

- Vocabulary: ["b", "h", "p", "#g", "#n", "#s", "#u", "#gs", "hu", "hug"]
- Corpus: ("hug", 10), ("p" "#u" "#g", 5), ("p" "#u" "#n", 12), ("b" "#u" "#n", 4), ("hu" "#gs", 5)
- **Test word : "hugs"**
- **Tokenization of "hugs" is ["hug", "#s"]**

Evaluating Summaries: ROUGE



ROUGE (Recall Oriented Understudy for Gisting Evaluation)

- Intrinsic metric for atomically evaluating summaries
- Based on BLEU (a metric used for machine translation)
- Not as good as human evaluation (“Did this answer the user’s question?”)
- But much more convenient

ROUGE-2

Given a document D, and an automatic summary X:

1. Have N humans produce a set of reference summaries of D
2. Run system, giving automatic summary X
3. What percentage of the bigrams from the reference summaries appear in X?

$$ROUGE-2 = \frac{\sum_{s \in \{\text{RefSummaries}\}} \sum_{\text{bigrams } i \in s} \min(\text{count}(i, X), \text{count}(i, S))}{\sum_{s \in \{\text{RefSummaries}\}} \sum_{\text{bigrams } i \in s} \text{count}(i, S)}$$

ROUGE-2 Example

Q: “What is water spinach?”

Human 1: Water spinach is a green leafy vegetable grown in the tropics.

Human 2: Water spinach is a semi-aquatic tropical plant grown as a vegetable.

Human 3: Water spinach is a commonly eaten leaf vegetable of Asia.

- System answer: Water spinach is a leaf vegetable commonly eaten in tropical areas of Asia.

$$\text{Rouge-2 score} = \frac{3 + 3 + 6}{10 + 9 + 9} = 12/28 = .43$$

Log-Likelihood Ratio (LLR)

innovate

achieve

lead

Centroid Based Summarization

1. Tokenize the input sentence S_j into set of non-stop words : W_i 's
2. Compute LLR for every W_i 's
3. Find the weights of every W_i 's
4. Using weight(W_i) compute weight(S) for every S_j 's
5. Rank the S_j 's in descending order of weight or compare with threshold θ
6. Pick top "k" sentences for summary

$\lambda(w) = \log$ likelihood ratio for a word is given by

$$\lambda = \frac{b(k, N, p)}{b(k_I, N_I, p_I).b(k_B, N_B, p_B)}$$

$$weight(w_i) = \begin{cases} 1 & \text{if } -2\log\lambda(w_i) > 10 \\ 1 & \text{if } w_i \in \text{question} \\ 0 & \text{otherwise} \end{cases}$$

$$weight(s) = \frac{1}{|S|} \sum_{w \in S} weight(w) \quad weight(s_i) = \sum_{w \in s_i} \frac{weight(w)}{|\{w | w \in s_i\}|}$$

S1 : Water spinach ipomoea aquatica semiaquatic leafy green plant long hollow stems spear heart shaped leaves = 8/15 = 0.53
S2 : grown Asia leaf vegetable = 2/4 = 0.5
S3 : leaves stems eaten fried flavored salt soups = 2/7 = 0.28
S4 : common names morning glory vegetable kangkong Malay = 3/7 = 0.42
S5 : related spinach closely related sweet potato convolvulus = 4/7 = 0.57

$$\begin{aligned} \text{Assume } P(W | I) &= 1/24, P(W | B) = 50/100000000 \\ P &= (1+50)/(24+100000000) \\ L(W) &= 0.000877 \rightarrow -2\log(L(w)) = 14.8 \end{aligned}$$

$$\begin{aligned} \text{Weight(vegetable)} &= 1 \\ \text{Weight}(S2) &= \frac{1}{4} * \text{sum_of_weights}\{\text{grown, asia, leaf, vegetable}\} \end{aligned}$$

Water spinach (ipomoea aquatica) is a semiaquatic leafy green plant with long hollow stems and spear or heart shaped leaves. It is not related to spinach, but is closely related to sweet potato and convolvulus.