

• Section 1: Linguistic Structure, Data and Analysis

Consider the following sentence: "the man with the hat and tie came after the man in the dark suit"

How many types and tokens are there in the sentence?

The number of tokens is the total number of words in the sentence. In this case, the number of tokens is 15.

The number of types is the total number of unique words in the sentence. In this case, the number of types is 12.

• Question 1B: Define a constituency (phrase) grammar and lexicon that analyses the following sentence by using the non-terminal symbols ' S, NP, VP, PP ' and the pre-terminal symbols 'Det, Noun, Verb, Prep'

The minister visited the power plant in the south of the country

Draw a constituency (phrase) structure tree for this sentence, using the grammar and lexicon you defined

The lexicon we need to analyse the sentence above is given by the types of the sentence. This is $L = \{\text{the}, \text{minister}, \text{visited}, \text{power}, \text{plant}, \text{in}, \text{south}, \text{of}, \text{country}\}$ together with their respectively non-terminal symbols:

~~the man~~

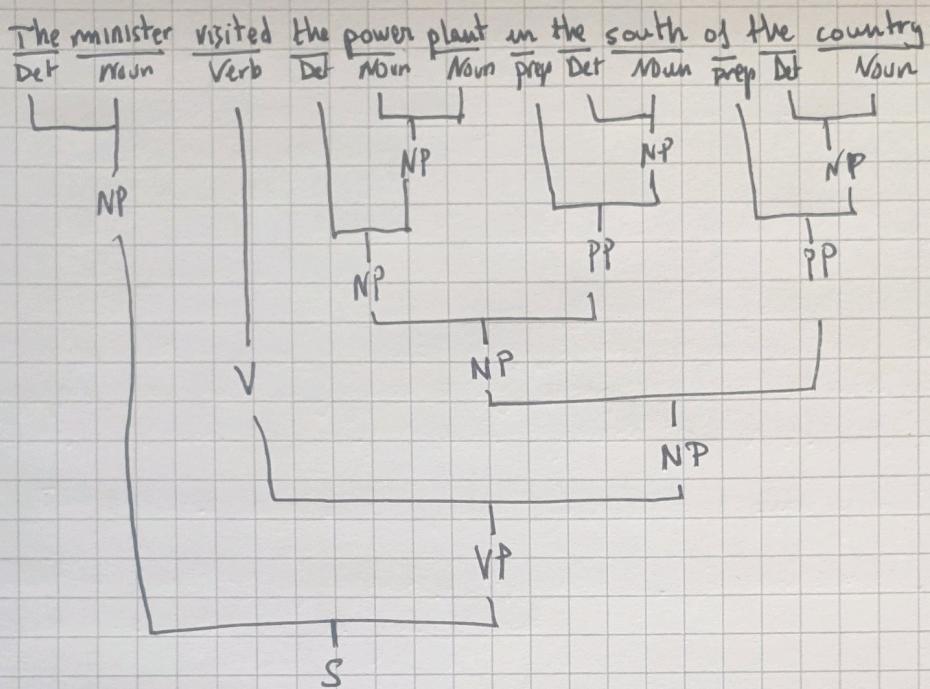
{
Det → the
Noun → minister, power, plant, south, country
Verb → visited
Prep → in, of

Now, with this set of terminal, non-terminal symbols and lexicon I define the following constituency grammar

{
NP → Det Noun | Noun Noun | Det NP | NP PP
PP → Prep Noun | Prep NP
VP → Verb | Verb NP | VPPP
S → NP VP

• Question 1B

With this grammar, the following constituency structure tree can be drawn from the given sentence:



• Question 1C: What is the difference between a parallel corpus and a comparable corpus?

A parallel corpus is a collection of translated documents while a comparable corpus is a collection of documents on the same topic in different languages

Section 2: Textual Similarity

s_1 : Tusk swipes at May for better border talk

s_2 : Tusk asks May for better border idea

Question 2A: Calculate the following similarities for s_1 and s_2 :

→ Dice similarity using a bag-of-words model

→ Jaccard similarity using a bag-of-words model

→ The length of the longest common subsequence

Jaccard similarity is given by $J_{s_1, s_2} = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$

where $|S_1 \cap S_2| = |\{\text{Tusk, May, for, better, border}\}| = 5$

$|S_1 \cup S_2| = |\{\text{Tusk, swipes, asks, at, May, for, better, border, talk, idea}\}| = 10$

Therefore $J_{s_1, s_2} = \frac{1}{2} = 0.5$

From here, it is easy to calculate the Dice similarity as $D = \frac{2S}{1+J} = \frac{1}{1.5} = \frac{2}{3} \approx 0.66$

Finally, the length of the longest common subsequence is 4 and is given by the length of "May for better border".

Question 2B: Recall the Damerau-Levenshtein Edit Distance:

$$d(i,j) = \min \begin{cases} d(i-1, j) + 1 \\ d(i, j-1) + 1 \\ d(i-1, j-1) + 1 & \text{if } S_{1,i} = S_{2,j} \\ d(i-2, j-2) + 1 & \text{if } S_{1,i} = S_{2,j-1} \wedge S_{1,i-1} = S_{2,j} \end{cases}$$

What is the Damerau-Levenshtein distance between these sentences?
Explain your method

Section 3: Language modelling

Consider the following poem by A.A. Milne as a corpus. Treat each line as a new sentence. Ignoring punctuation, it is 100 words long.

The wind on the hill.

No one can tell me nobody knows where the wind comes from where the wind goes.

It's flying from somewhere as fast as I can I couldn't keep up with it not if I ran.

But if I stopped holding the string of my kite, it would blow with the wind for a day and a night.

And then when I found it, wherever it blew, I should know that the wind had been going there too.

So then I could tell them Where the wind goes

But where the wind comes from nobody knows

Question 3A: Calculate the unigram probabilities ignoring case for the words: "been", "had", "the", "wind", "where".

As there is a total of 100 words, given a word w , the unigram probability of w can be calculated by:

$$P(w) = \frac{\text{#number of occurrences of } w}{100}$$

This is:

$$P(\text{"been"}) = \frac{2}{100} \quad P(\text{"wind"}) = \frac{7}{100}$$

$$P(\text{"had"}) = \frac{2}{100} \quad P(\text{"where"}) = \frac{4}{100}$$

$$P(\text{"the"}) = \frac{8}{100}$$

Question 3B: Calculate the bigram probability ignoring case for the combinations that are not provided in the following table.

$p(w_2 w_1)$	$w_1 = \text{been}$	$w_1 = \text{had}$	$w_1 = \text{the}$	$w_1 = \text{where}$	$w_1 = \text{wind}$
$w_1 = \text{been}$	0	0	0	0	0
$w_1 = \text{had}$	1	0	0	0	0
$w_1 = \text{the}$	0	0	0	0	$\frac{7}{8}$
$w_1 = \text{where}$	0	0	1	0	0
$w_1 = \text{wind}$	0	$\frac{1}{7}$	0	0	0

• Question 3B

The probability of a bigram (w_1, w_2) can be calculated as :

$$P(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\text{count}(w_1)}$$

$$\begin{aligned} \text{Count(had, been)} &= 1, & \text{count(had)} &= 1 \Rightarrow P(\text{been}|\text{had}) = 1 \\ \text{Count(the, where)} &= 0, & \text{count(the)} &= 8 \Rightarrow P(\text{where}|\text{the}) = 0 \\ \text{Count(the, wind)} &= 7, & \text{count(the)} &= 8 \Rightarrow P(\text{wind}|\text{the}) = 7/8 \\ \text{Count(where, the)} &= 4, & \text{count(where)} &= 4 \Rightarrow P(\text{where}|\text{the}) = 1 \\ \text{Count(wind, had)} &= 1, & \text{count(wind)} &= 7 \Rightarrow P(\text{had}|\text{wind}) = 1/7 \end{aligned}$$

• Question 3C:

State the formula for a bigram language model applied to the sentence "the wind had been". Using this bigram language model calculate the probability of the line "the wind had been".

In general, using a bigram language model we can approximate the probability of a sequence of words by :

$$P(w_1 \dots w_n) \approx p(w_n|w_{n-1}) \cdot p(w_{n-1}|w_{n-2}) \cdots p(w_2|w_1) p(w_1)$$

$$\text{In this specific sentence, this is } p(\text{"the wind had been"}) = \underbrace{p(\text{been}|\text{had})}_{1} \underbrace{p(\text{had}|\text{wind})}_{1/7} \underbrace{p(\text{wind}|\text{the})}_{7/8} \underbrace{p(\text{the})}_{?/8}$$

$$\text{This is } p(\text{"the wind had been"}) = \frac{1}{8} \cdot \frac{1}{7} \cdot \frac{7}{8} \text{ where } p(\text{"the"}) = p(\text{"the"}|\text{start}) = \frac{1}{7}$$

$$\text{Therefore } p(s) = \frac{1}{56}$$

• Question 3D:

$p(\text{"The wind had been there"}) = 0$ given the bigram model. Briefly explain why and suggest a model that produces a non-zero probability for this sentence.

$$P(\text{"The wind had been there"}) = P(\text{"there"}|\text{"been"}) P(\text{"been"}|\text{"had"}) P(\text{"had"}|\text{"wind"}) P(\text{"wind"}|\text{"the"}) P(\text{"the"}|\text{"start"})$$

As the bigram ("been", "there") does not exist in the corpus this probability will be always zero ($P(\text{"there"}|\text{"been"}) = 0$).

An alternative would be using Add-One Smoothing on which any bigram is counted one more time than the occurrences in the corpus.

In this case

$$P(\text{"there"}|\text{"been"}) = \frac{c(\text{been}, \text{there}) + 1}{r(\text{been}) + \frac{\# \text{bigrams}}{\# 1\text{-grams}}} = \frac{1}{1 + 99} = \frac{1}{99}$$

- Section 4: Information Extraction

Consider the following sentences

S₁ = Shares in Smurfit Kappa have risen by over 18 pc in early trading on the London Stock Exchange

S₂ = Adidas shares were up 1.6%, marking the biggest increase among the largest shares in Germany.

S₃ = Shares of Jaypee Infratech climbed over 4 per cent on Tuesday morning.

S₄ = United Technologies shares rise 2% as a "well-known" activist takes a position

S₅ = "Shares of Commonwealth Bank declined 0.92 percent"

- Question 4A: Annotate sentences S₂ and S₄, with Named Entities of type COMPANY, NUMBER, COUNTRY and TIME where appropriate. Use the following annotation format:

[COMPANY Adidas] shares were up [NUMBER 1.6]%. - ..

[COMPANY Adidas] shares were up [NUMBER 1.6]%, marking the biggest increase among the largest shares in [COUNTRY Germany]

[Shares of] [COMPANY Jaypee Infratech] climbed over [NUMBER 4] per cent on [TIME Tuesday morning].

- Question 4B: Assume we want to extract the following relations from the sentences above:

Shares-Up-Percentage([COMPANY, NUMBER])

Shares-Down-Percentage([COMPANY, NUMBER])

For instance Shares-Up-Percentage(Smurfit Kappa, 18)

Provide patterns that can be used to extract this information for all companies mentioned in the sentences above

Hearst Patterns for Shares-Up-Percentage:

[COM] have risen over [NUM], [COM] were up [NUM], [COM] shares rise [NUM],
[COM] climbed over [NUM]

for Shares-Down-Percentage: [COM] declined [NUM]

? Question 4C What is the core difference between 'open information extraction' (Open IE) and 'knowledge base population (KBP)'?

Open IE refers to the extraction of relation tuples, typically binary relations, from plain text (Unsupervised Relation Extraction)

KBP refers to the task of discovering new facts about entities from a large corpus, and augmenting a knowledge base with these facts (Supervised / based on a knowledge graph)