

Start

End

Send

learning

game

Sol 1(a.) Tokenization is the process of breaking down a text into smaller units called tokens.

These tokens can be words, characters or subwords depending on the chosen tokenization scheme.

Tokenization is a fundamental step in NLP tasks, as it allows the model to process and understand the text at a more granular level.

BPE (Byte-Pair Encoding) is a popular tokenization algorithm that uses a statistical approach to build a vocabulary of subword units.

It works by iteratively merging the most frequent pairs of characters or character sequences in a corpus.

Advantages of BPE :-

- Subword representation for handling OOV words
- Improved generalization by recognizing shared subwords
- Vocabulary compression for reduced memory usage.
- Captures morphological information in languages.

Training the BPE algo. using the below mentioned table:

dictionary

low → 5 → 3 → 2 → 1 → 0

lowest → 2 → 3 → 2 → 1 → 0

newer → 6 → 5 → 4 → 3 → 2 → 1 → 0

wider → 3 → 2 → 1 → 0

new → 2 → 1 → 0

Vocabulary (11 letters)

- , d, e, i, l, n, o, s, t, w

Note:- Each word is represented as a sequence of characters plus a special end-of-word symbol here i.e., — (underline)

Correct order of iterations are :-

(remember, in each iteration)
(we have to count pairs from dictionary again.)

Initially

dictionary

vocabulary (11 letters)

5	low -	b, l, d, e, i, s, l, n, o, r, s, t
2	low e st -	, w,
6	new e r -	
3	wid e r -	
2	new -	

⇒ we first count all pairs of symbols - the most freq. is r -
(9 times)
merge r & s - treating rs as one symbol

next iteration

dictionary

vocabulary

5	l o w -	, d, e, i, s, l, n, o, r, s, t
2	l o w e s t -	, t, w, r,
6	n e w e r -	
3	w i d e r -	
2	n e w -	

Now, most frequent pair is ew - i.e., make them ex-

next iteration

dictionary

vocabulary

5	l o w -	, d, e, i, s, l, n, o, r, s, t
2	l o w e s t -	, t, w, r, ex-
6	n e w ex -	
3	w i d e ex -	
2	n e w -	

⇒ Now, mfp is ew.

~~next iteration~~

Dictionary

5 l o w -
 2 l o w e s t -
 6 n ew ex-
 3 w i d ex-
 2 n ew -

Vocabulary

-d, e, i, l, n, o, r,
 s, t, w, x-, ex-, ew
 w
 x o w e m
 x o w e m
 x o w e m

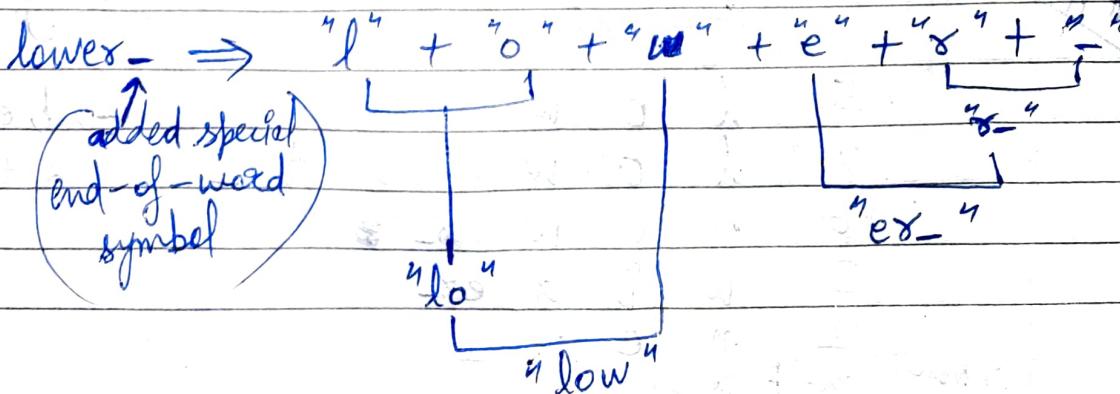
⇒ now, mfp is new.

If we continue, the next merges are :-

Merge	Current Vocabulary
(n, ew)	-d, e, i, l, n, o, r, s, t, w, x-, ex-, ew, new
(l, o)	l, o,
(l, o, w)	low
(new, ex-)	newer-
(low, -)	low-

When we need to tokenize a test sentence, we just run the merges we have learned, greedily, in the order we learned them, on the test data.
 (Thus, the freq.'s in test data don't play a role.)
 Just, freq.'s in training data are vital.

Now, Testing :- the final tokenization of the word 'lower'



hey, however :- "law" + "er-"

so, A new (unknown) testing word "lomer-" would be merged into two tokens i.e., "law" + "er-".

Sol 1(b.) Bigram counts :- (given)

	<u>i</u>	<u>want</u>	<u>to</u>	<u>eat</u>	<u>chinese</u>	<u>food</u>	<u>lunch</u>	<u>spend</u>
<u>i</u>	5	827	0	9	0	0	0	23
<u>want</u>	2	0	608	1	6	6	5	1
<u>to</u>	2	0	4	686	2	0	6	211
<u>eat</u>	0	0	2	0	16	2	42	0
<u>chinese</u>	1	0	0	0	0	82	1	0
<u>food</u>	15	0	15	0	0	4	0	0
<u>lunch</u>	2	0	0	0	0	1	0	0
<u>spend</u>	1	0	1	0	0	0	0	0

I am assuming the whole corpus has V unique words.
i.e., $V = 1446$ out of which only 8 words bigram counts are given to us.

Unigram counts :- (assumed)

<u>i</u>	<u>want</u>	<u>to</u>	<u>eat</u>	<u>chinese</u>	<u>food</u>	<u>lunch</u>	<u>spend</u>
2533	927	2417	746	158	1093	341	278

$$P(w_m | w_{m-1}) = \frac{C(w_{m-1} w_m)}{C(w_{m-1})}$$

$$P(i|i) = \frac{C(ii)}{C(i)} = \frac{5}{2533} = 0.00197 \\ \approx 0.002$$

$$P(want | i) = \frac{C(want i)}{C(want)} = \frac{2}{927} = 0.002157 \\ \approx 0.0022$$

similarly for others.

Bigram Probability matrix

i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33					
want	0.0022	0					
to	0.00083	0					
eat	0	0					
chinese		0					
food		0					
lunch		0					
spend		0					

Similarly, you can make table, for Bigram Probability (matrix)

with add-1 smoothing, using below mentioned formula.

$$P^*(w_m | w_{m-1}) = \frac{C(w_{m-1} w_m) + 1}{C(w_{m-1}) + V}$$

Laplace

Sol 1(c.)

(1) Segmentation :-

It is the process of dividing a sentence or a document into smaller units like words or phrases or sentences for further analysis.

This is a common technique used in natural language processing and text analysis.

Ex:-

"I love chatting with friends"
⇒ ["I", "love", "chatting", "with", "friends"].

(2) Lemmatization :-

It is the process of reducing words to their base or root or dictionary form known as the lemma.

This is done to simplify the analysis of text by grouping ~~those~~ together words that have the same base form.

Ex:- "am", "is", "are", "was", "were" gets reduced in Lemmatization to their base form "be".

(3) Stemming :-

It is the process of reducing words to their root form by removing the suffixes.

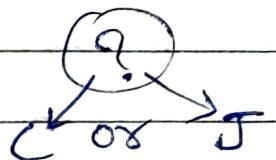
This is a simpler approach than lemmatization because it does not consider the context of the word.

Ex:- "Jumping" gets reduced to "Jump".

Sol 2(a.)

	<u>Doc(d)</u>	<u>Words</u>	<u>Class</u>
Training	1	Chinese, Beijing, Chinese	C
	2	Chinese, Chinese, Shanghai	C
	3	Chinese, Macao	C
	4	Tokyo, Japan, Chinese	J
Test	5	Chinese, Chinese, Tokyo , Tokyo, Japan	?

To compute :- most likely class for Doc 5.



Assume, a multinomial naive Bayes classifier and use add- α Laplace smoothing for the likelihoods. say $\alpha = 1$.

a class

$$P(\vec{X}) = \frac{\text{Count of } \vec{X}}{N}, N = \text{total count of classes.}$$

total no. of words in the c class.

$$P(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$$

$\uparrow \quad \uparrow$

total no. of words
in c class

vocabulary size
(size of unique words dictionary)
(size of unique words from training dataset)

Priors :-

$$P(C) = \frac{3}{4} \quad (\text{Probability of class C})$$

$$P(J) = \frac{1}{4} \quad (\text{--- J})$$

Conditional Probabilities :-

$$P(\text{Chinese} | C) = \frac{(5+1)}{(8+6)} = \frac{6}{14}$$

$$P(\text{Tokyo} | C) = \frac{(0+1)}{(8+6)} = \frac{1}{14}$$

$$P(\text{Japan} | C) = \frac{(0+1)}{(8+6)} = \frac{1}{14}$$

$$P(\text{Chinese} | J) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

$$P(\text{Tokyo} | J) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

$$P(\text{Japan} | J) = \frac{(1+1)}{(3+6)} = \frac{2}{9}$$

Predicting class for Doc 5.

Probability of Docs to come in C class

$$P(C | ds) = (P(C))! \cdot (P(\text{Chinese} | C))^{\frac{2}{9}} \cdot (P(\text{Tokyo} | C))^{\frac{2}{9}} \cdot (P(\text{Japan} | C))^{\frac{2}{9}}$$

$$= \frac{3}{4} \cdot \left(\frac{6}{14}\right)^2 \cdot \frac{1}{14} \cdot \frac{1}{14} = 0.0007028$$

Probability of Docs to come in J class

$$P(J | ds) = \frac{1}{4} \cdot \left(\frac{2}{9}\right)^2 \cdot \frac{2}{9} \cdot \frac{2}{9} = 0.0006097$$

$$\therefore P(C | ds) > P(J | ds)$$

∴ Doc5 will come into class C.

Sol 2(b.)Training data

<u>S. No.</u>	<u>Document</u>	<u>Class</u>
1	Natural Language Processing	A
2	Language Model Learning	A
3	Ngram Language Model	A
4	Text Classification Model	A
5	Text Processing Model	A
6	Computer Vision	B
7	Image Classification Model	B
8	Object Segmentation	B
9	Image Processing	B
10	Object Recognition	B

Testing data

<u>S. No.</u>	<u>Document</u>	<u>Class</u>
1	Object Recognition Model	?
2	Text Recognition Model	?

Using Multinomial Naive Bayes (Unigram Model)

$$P(A) = \frac{5}{10} = \frac{1}{2}, \quad P(B) = \frac{5}{10} = \frac{1}{2}, \quad |V| = \text{Unique Words} = 14$$

Now, test data :- Object Recognition Model.

$$P(\text{object}|A) = \frac{0+1}{15+14} = \frac{1}{29}$$

$$P(\text{object}|B) = \frac{2+1}{11+14} = \frac{3}{25}$$

$$P(\text{recognition}|A) = \frac{0+1}{15+14} = \frac{1}{29}$$

$$P(\text{recognition}|B) = \frac{1+1}{11+14} = \frac{2}{25}$$

$$P(\text{Model}|A) = \frac{4+1}{15+14} = \frac{5}{29}$$

$$P(\text{Model}|B) = \frac{1+1}{11+14} = \frac{2}{25}$$

(d)

Predicting class for Test doc 1.

$$P(A|d_1) = (P(A))' \cdot (P(\text{object}|A))' \cdot (P(\text{recognize}|A))' \cdot (P(\text{Model}|A))'$$

$$= \frac{1}{2} \cdot \frac{1}{29} \cdot \frac{1}{29} \cdot \frac{5}{29} = 0.0001$$

$$P(B|d_1) = (P(B))' \cdot (P(\text{object}|B))' \cdot (P(\text{recognize}|B))' \cdot (P(\text{Model}|B))'$$

$$= \frac{1}{2} \cdot \frac{3}{25} \cdot \frac{2}{25} \cdot \frac{2}{25} = 0.0003$$

∴ d₁ :- Object Recognition Model will come into class 'B'

✓

Now test data :- Text Recognition Model

$$P(\text{Text}|A) = \frac{2+1}{15+14} = 3/29 \quad P(\text{Text}|B) = \frac{2+1}{11+14} = 1/25$$

$$P(\text{Recognize}|A) = \frac{0+1}{15+14} = 1/29 \quad P(\text{Recognize}|B) = \frac{1+1}{11+14} = 2/25$$

$$P(\text{Model}|A) = \frac{4+1}{15+14} = 5/29 \quad P(\text{Model}|B) = \frac{2}{11+14} = 2/25$$

Predicting class for Test doc (d) 2

$$P(A|d_2) = (P(A))' \cdot (P(\text{Text}|A))' \cdot (P(\text{recognize}|A))' \cdot (P(\text{Model}|A))'$$

$$= \frac{1}{2} \cdot \frac{3}{29} \cdot \frac{1}{29} \cdot \frac{5}{29} = 0.0003$$

$$P(B|d_2) = \frac{1}{2} \cdot \frac{1}{25} \cdot \frac{2}{25} \cdot \frac{2}{25} = 0.0001$$

∴ d₂ :- Text Recognition Model will come into class 'A'.

Using Multinomial Naive Bayes (Bigram Model)

Aniket

Date _____
Page _____

<2> Natural Language Processing </2>

A

<2> Language Model Learning </2>

A

<2> Ngram Language Model </2>

A

<2> text classification Model </2>

A

<2> Text Processing Model </2>

A

<2> Computer Vision </2>

B

<2> Image Classification Model </2>

(sh) B

<2> Object Segmentation </2>

B

<2> Image Processing </2>

B

<2> Object Recognition </2>

B

If we include <2> & <3>, we get

~~we have 14 words including <2> & if we include <3>~~
~~(this can also appear as second element of a bigram),~~
we get $|V| = 16$ for our vocabulary.

$$(14 + 1) + 1 \\ (2>) \quad (3>)$$

Test doc 1 :- <2> object recognition Model </2>

$$\begin{aligned} P(A|d_1) &= P(A) \cdot P(\text{object} | <2>) \cdot P(\text{recognition} | \text{object}) \\ &\quad \cdot P(\text{Model} | \text{recognition}) \cdot P(</2> | \text{Model}) \\ &= P(A) \cdot \frac{\text{Count}(<2>, \text{object}) + 1}{\text{Count}(<2>) + |V|} \quad \dots \text{(similarly)} \\ &= \frac{5}{10} \cdot \left(\frac{0+1}{5+16} \right) \cdot \left(\frac{0+1}{0+16} \right) \cdot \left(\frac{0+1}{0+16} \right) \cdot \left(\frac{3+1}{4+16} \right) = 0.000019 \end{aligned}$$

$$P(B|d1) = \frac{5}{10} \cdot \left(\frac{2+1}{5+16}\right) \cdot \left(\frac{1+1}{2+16}\right) \cdot \left(\frac{0+1}{1+16}\right) \cdot \left(\frac{1+1}{1+16}\right) = 0.000055$$

$\therefore d1$:- "Object Recognition Model" will come ^{into} ~~into~~ class 'B'.

Test doc2 :- <8> Text recognition model </8>

$$P(A|d2) = \frac{5}{10} \cdot \left(\frac{2+1}{5+16}\right) \cdot \left(\frac{0+1}{2+16}\right) \cdot \left(\frac{0+1}{0+16}\right) \cdot \left(\frac{3+1}{4+16}\right) = 0.00005$$

$$P(B|d2) = \frac{5}{10} \cdot \left(\frac{0+1}{5+16}\right) \cdot \left(\frac{0+1}{0+16}\right) \cdot \left(\frac{0+1}{1+16}\right) \cdot \left(\frac{1+1}{1+16}\right) = 0.00001$$

$\therefore d2$:- "Text Recognition Model" will come ^{into} class "A".