

Name: Muhammad Ghazanfar

Roll#: 20P-0567

Sec: 7A

Course: NLP

Lab#02

## Task 1: Bigrams & Trigrams

```
Task#01
                                                                                                                                                                                                                                       Python
            from nltk.book import *
 \cdots *** Introductory Examples for the NLTK Book ***
       Loading text1, ..., text9 and sent1, ..., sent9
Type the name of the text or sentence to view it.
Type: 'texts()' or 'sents()' to list the materials.
text1: Moby Dick by Herman Melville 1851
text2: Sense and Sensibility by Jane Austen 1811
        text4: Inaugural Address Corpus
text5: Chat Corpus
        text6: Monty Python and the Holy Grail
text7: Wall Street Journal
text8: Personals Corpus
        text9: The Man Who Was Thursday by G . K . Chesterton 1908
             words = sorted(set(text1))[280:]
           longwords = [w for w in words if len(w) > 16]
print(longwords)
           Magh_freq = [w for w in words if fdistl[w] > 500]
print(high_freq)
          eign_words = [w for w in words if w.endswith('eign')]
print(eign_words)
                                                                                                                                                                                                                                       Python
           for w in words:
    if w.endswith('eign'):
        print(w)
       sovereian
        list(bigrams(text1))[:10]
                                                                                                                                                                                                                                       Python
··· [('[', 'Moby'),
('Moby', 'Dick'),
('Dick', 'by'),
```

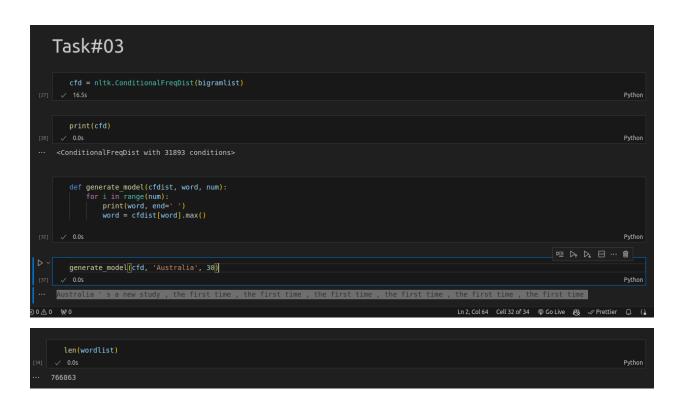
```
Sperm Whale; Moby Dick; White Whale; old man; Captain Ahab; sperm whale; Right Whale; Captain Peleg; New Bedford; Cape Horn; cried Ahab; years ago; lower jaw; never mind; Father Mapple; cried Stubb; chief mate; white whale; ivory leg; one hand
                         from nltk.util import ngrams
                        list(ngrams(text1, 3))
   ...
[('[', 'Moby', 'Dick'),
   ('Moby', 'Dick', 'by'),
   ('Dick', 'by', 'Herman'),
   ('by', 'Herman', 'Melville'),
   ('Herman', 'Melville', '1851'),
   ('Melville', '1851', ']'),
   ('1851', ']', 'ETYMOLOGY'),
      ... [('olive', 'leaf', 'pluckt'),
    ('sewed', 'fig', 'leaves'),
    ('yield', 'royal', 'dainties'),
    ('Fifteen', 'cubits', 'upward'),
    ('leaf', 'pluckt', 'o')]
⊗ 2 ∆ 3 🙀 0
                                                                                                                                                                                                                                                                                                                                                                    Cell 2 of 45 @ Go Live 🗞 🚨 🛟
```

	Text1	Text2	Text3	
10 frequently occuring Bigrams	Sperm Whale; Moby Dick; White Whale; old man; Captain Ahab; sperm whale; Right Whale; Captain Peleg; New Bedford; Cape Horn;	Colonel Brandon; Sir John; Lady Middleton; Miss Dashwood; every thing; thousand pounds; dare say; Miss Steeles; said Elinor; Miss Steele;	said unto; pray thee; thou shalt; thou hast; thy seed; years old; spake unto; thou art; LORD God; every living;	
5 frequently occuring Trigrams	[('AFTER', 'EXCHANGING', 'HAILS'), ('Anacharsis', 'Clootz', 'deputation'), ('CAULKING', 'ITS', 'SEAMS'), ('ELIZABETH', 'OAKES', 'SMITH'), ('Et', 'tu', 'Brute')]	[('Austen', '1811', ']'), ('Jane', 'Austen', '1811'), ('200', 'L', 'per'), ('Drury', 'Lane', 'lobby'), ('L', 'per', 'annum')]	[('olive', 'leaf', 'pluckt'), ('sewed', 'fig', 'leaves'), ('yield', 'royal', 'dainties'), ('Fifteen', 'cubits', 'upward'), ('leaf', 'pluckt', 'o')]	
Number of words with length > 16	['cannibalistically',    'characteristically',    'circumnavigations',    'comprehensiveness',    'indispensableness',    'preternaturalness',    'subterraneousness',    'superstitiousness',    'uncomfortableness',    'uncompromisedness',    'uninterpenetratingly']	['companionableness', 'disinterestedness', 'disqualifications']	[] NO words greater than len 16	
Number of words with frequency > 500	['Ahab', 'But', 'I', 'The', 'a', 'all', 'an', 'and', 'are', 'as', 'at', 'be', 'but', 'by', 'for', 'from', 'had', 'have', 'he', 'him', 'his', 'in', 'into', 'is', 'it', 'like', 'man', 'me', 'more',	["'", ',', '.', ';', 'I', 'a', 'all', 'and', 'as', 'at', 'be', 'but', 'by', 'could', 'for', 'from', 'had', 'have', 'he', 'her', 'him',	[',', '.', ';', 'A nd', 'and', 'he', 'his', 'in', 'of', 'that', 'the', 'to', 'unto']	

	<pre>'my', 'not', 'now', 'of', 'on', 'one', 'or', 'out', 's', 'ship', 'so', 'some', 'that', 'the', 'their', 'then', 'there', 'they', 'this', 'to', 'up', 'upon', 'was', 'were', 'whale', 'when', 'which', 'with', 'you']</pre>	'his', 'in', 'is', 'it', 'my', 'not', 'of', 'on', 's', 'she', 'so', 'that', 'the', 'to', 'was', 'which', 'with', 'would', 'you']	
Number of words ending in "ed"	2196	902	281

## Task 2: Accessing Corpora

## Task 3: Generating Random Text with Bigrams



Number of Words In Corpus	30 word generated sentence
766863	Today, the first time,
766863	Australia 's a new study , the first time