Lab 02 NLTK: Language Models Task 1: Bigrams & Trigrams In [3]: import nltk from nltk.book import \* \*\*\* 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 text3: The Book of Genesis 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 In [4]: words = sorted(set(text1))[280:] In [5]: print(len(words)) 19037 In [6]: longwords = [w for w in words if len(w) > 16] In [7]: longwords ['cannibalistically', 'characteristically', 'circumnavigations', 'comprehensiveness', 'indispensableness', 'preternaturalness', 'subterraneousness', 'superstitiousness', 'uncomfortableness', 'uncompromisedness', 'uninterpenetratingly'] In [8]: fdist1 = FreqDist(text1) In [9]: high\_freq = [w for w in words if fdist1[w] > 500] In [10]: print(high\_freq) ['Ahab', 'But', 'I', 'The', 'a', 'all', 'an', 'and', 'are', 'as', 'at', 'be', 'but', 'by', 'for', 'from', 'had', 'have', 'he', 'in', 'into', 'is', 'it', 'like', 'man', 'me', 'more', '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'] In [12]: **import** latexify In [13]: **import** math @latexify.function def IDF(corpus\_list, token): document\_frequency = sum(1 for corpus in corpus\_list if token in corpus) # Add 1 to both numerator and denominator for smoothing inverse\_document\_frequency = math.log((len(corpus\_list) + 1) / (document\_frequency + 1), 10) return inverse\_document\_frequency In [14]: IDF Out[14]:  $document\_frequency = \sum_{(corpus \in corpus\_list) \land (token \in corpus)} (1)$  $ext{inverse\_document\_frequency} = \log\Bigl(rac{ ext{len(corpus\_list)} + 1}{ ext{document\_frequency} + 1}, 10\Bigr)$  $IDF(corpus\_list, token) = inverse\_document\_frequency$ In [15]: high\_freq\_idf = [w for w in words if IDF(text1, w) > 1] In [16]: len(high\_freq\_idf) Out[16]: 19027 In [17]: print(high\_freq\_idf[:10]) ['ABOUT', 'ACCOUNT', 'ADDITIONAL', 'ADVANCING', 'ADVENTURES', 'AFFGHANISTAN', 'AFRICA', 'AFTER', 'AGAINST', 'AHAB'] In [18]: eign\_words = [w for w in words if w.endswith('eign') ] In [19]: eign\_words ['Sovereign', 'foreign', 'reign', 'sovereign'] Out[19]: In [20]: **for** w **in** words: if w.endswith('eign'): print(w) Sovereign foreign reign sovereign In [21]: len(list(bigrams(text1))) Out[21]: In [22]: text1.collocation\_list() [('Sperm', 'Whale'), Out[22]: ('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')] In [23]: **from** nltk.util **import** ngrams list(ngrams(text1, 3))[:10] Out[23]: [('[', 'Moby', 'Dick'), ('Moby', 'Dick', 'by'), ('Dick', 'by', 'Herman'), ('by', 'Herman', 'Melville'), ('Herman', 'Melville', '1851'), ('Melville', '1851', ']'), ('1851', ']', 'ETYMOLOGY'), (']', 'ETYMOLOGY', '.'), ('ETYMOLOGY', '.', '('), ('.', '(', 'Supplied')] In [24]: **from** nltk.collocations **import** TrigramCollocationFinder, TrigramAssocMeasures TrigramCollocationFinder.from\_words(text1).nbest(TrigramAssocMeasures().pmi, 10) [('AFTER', 'EXCHANGING', 'HAILS'), ('Anacharsis', 'Clootz', 'deputation'), ('CAULKING', 'ITS', 'SEAMS'), ('ELIZABETH', 'OAKES', 'SMITH'), ('Et', 'tu', 'Brute'), ('Ex', 'officio', 'professors'), ('Fogo', 'Von', 'Slack'), ('Ganders', 'formally', 'indite'), ('Kentucky', 'Mammoth', 'Cave'), ('LANTERNS', 'BUSILY', 'FILING')] 10 frequently occuring Bigrams In [25]: **from** nltk **import** bigrams, FreqDist def top\_frequent\_bigrams(text, n = 10): bigrams\_list = list(bigrams(text)) bigram\_freq = FreqDist(bigrams\_list) return bigram\_freq.most\_common(n) text1\_bigrams = top\_frequent\_bigrams(text1) text2\_bigrams = top\_frequent\_bigrams(text2) text3\_bigrams = top\_frequent\_bigrams(text3) print("Top 10 Frequently Occurring Bigrams in Text1:") print(text1\_bigrams) print("\nTop 10 Frequently Occurring Bigrams in Text2:") print(text2\_bigrams) print("\nTop 10 Frequently Occurring Bigrams in Text3:") print(text3\_bigrams) Top 10 Frequently Occurring Bigrams in Text1: [((',', 'and'), 2607), (('of', 'the'), 1847), (("'", 's'), 1737), (('in', 'the'), 120), ((',', 'the'), 908), ((';', 'and'), 853), (('to', 'the'), 712), (('.', 'But'), 596), ((',', 'that'), 584), (('.', '"'), 557)] Top 10 Frequently Occurring Bigrams in Text2: [((',', 'and'), 1598), (("'", 's'), 700), ((';', 'and'), 605), (('Mrs', '.'), 529), (('of', 'the'), 430), (('to', 'be'), 428), (('."', '""), 392), (('.', '""), 369), (('in', 'the'), 348)] Top 10 Frequently Occurring Bigrams in Text3: [((',', 'and'), 1491), (('.', 'And'), 1038), (('of', 'the'), 372), (('in', 'the'), 287), ((';', 'and'), 262), (('said', ','), 259), (("'", 's'), 255), (('And', 'he'), 192), (('And', 'the'), 185), (('said', 'unto'), 178)] 5 frequently occuring Trigrams In [26]: **from** nltk **import** trigrams, FreqDist from nltk.book import text1, text2, text3 def top\_frequent\_trigrams(text, n=5): trigrams\_list = list(trigrams(text)) trigram\_freq = FreqDist(trigrams\_list) return trigram\_freq.most\_common(n) text1\_trigrams = top\_frequent\_trigrams(text1, n=5) text2\_trigrams = top\_frequent\_trigrams(text2, n=5) text3\_trigrams = top\_frequent\_trigrams(text3, n=5) print("\nTop 5 Frequently Occurring Trigrams in Text1:") print(text1\_trigrams) print("\nTop 5 Frequently Occurring Trigrams in Text2:") print(text2\_trigrams) print("\nTop 5 Frequently Occurring Trigrams in Text3:") print(text3\_trigrams) Top 5 Frequently Occurring Trigrams in Text1: [((',', 'and', 'the'), 187), (('don', "'", 't'), 103), (('of', 'the', 'whale'), 101), ((',', 'in', 'the'), 93), ((',', 'then', ','), 87)] Top 5 Frequently Occurring Trigrams in Text2: [(('Mrs', '.', 'Jennings'), 230), (('Mrs', '.', 'Dashwood'), 121), ((',', 'however', ','), 88), ((',', 'and', 'the'), 87), (('.', 'Mrs', '.'), 80)] Top 5 Frequently Occurring Trigrams in Text3: [(('.', 'And', 'he'), 162), (('.', 'And', 'the'), 158), (('the', 'land', 'of'), 101), (('he', 'said', ','), 86), (('And', 'he', 'said'), 84)] Number of Words with Length > 16: In [27]: def count\_words\_length\_greater\_than(text, length=16): long\_words = [word for word in text if len(word) > length] return len(long\_words) text1\_long\_words\_count = count\_words\_length\_greater\_than(text1) text2\_long\_words\_count = count\_words\_length\_greater\_than(text2) text3\_long\_words\_count = count\_words\_length\_greater\_than(text3) print("\nNumber of Words with Length > 16 in Text1:", text1\_long\_words\_count) print("Number of Words with Length > 16 in Text2:", text2\_long\_words\_count) print("Number of Words with Length > 16 in Text3:", text3\_long\_words\_count) Number of Words with Length > 16 in Text1: 14 Number of Words with Length > 16 in Text2: 3 Number of Words with Length > 16 in Text3: 0 Number of words with frequency > 500 In [28]: def count\_words\_frequency\_greater\_than(text, frequency=500): word\_freq = FreqDist(text) high\_frequency\_words = [word for word, freq in word\_freq.items() if freq > frequency] return len(high\_frequency\_words) text1\_high\_freq\_words\_count = count\_words\_frequency\_greater\_than(text1) text2\_high\_freq\_words\_count = count\_words\_frequency\_greater\_than(text2) text3\_high\_freq\_words\_count = count\_words\_frequency\_greater\_than(text3) print("\nNumber of Words with Frequency > 500 in Text1:", text1\_high\_freq\_words\_count) print("Number of Words with Frequency > 500 in Text2:", text2\_high\_freq\_words\_count) print("Number of Words with Frequency > 500 in Text3:", text3\_high\_freq\_words\_count) Number of Words with Frequency > 500 in Text1: 67 Number of Words with Frequency > 500 in Text2: 45 Number of Words with Frequency > 500 in Text3: 13 Number of ending in "ed" In [29]: def count\_words\_ending\_in\_ed(text): ed\_words = [word for word in text if word.endswith("ed")] return len(ed\_words) text1\_ed\_words\_count = count\_words\_ending\_in\_ed(text1) text2\_ed\_words\_count = count\_words\_ending\_in\_ed(text2) text3\_ed\_words\_count = count\_words\_ending\_in\_ed(text3) print("\nNumber of Words Ending in 'ed' in Text1:", text1\_ed\_words\_count) print("Number of Words Ending in 'ed' in Text2:", text2\_ed\_words\_count) print("Number of Words Ending in 'ed' in Text3:", text3\_ed\_words\_count) Number of Words Ending in 'ed' in Text1: 8491 Number of Words Ending in 'ed' in Text2: 4866 Number of Words Ending in 'ed' in Text3: 1238 Task 2: Accessing Corpora nltk.corpus.gutenberg.fileids() ['austen-emma.txt', 'austen-persuasion.txt', 'austen-sense.txt', 'bible-kjv.txt', 'blake-poems.txt', 'bryant-stories.txt', 'burgess-busterbrown.txt', 'carroll-alice.txt', 'chesterton-ball.txt', 'chesterton-brown.txt', 'chesterton-thursday.txt', 'edgeworth-parents.txt', 'melville-moby\_dick.txt', 'milton-paradise.txt', 'shakespeare-caesar.txt', 'shakespeare-hamlet.txt', 'shakespeare-macbeth.txt', 'whitman-leaves.txt'] In [31]: gutenberg\_sc = nltk.corpus.gutenberg.words('shakespeare-caesar.txt') In [32]: len(gutenberg\_sc) Out[32]: 25833 In [33]: **from** nltk.corpus **import** brown brown.words() ['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', ...] In [34]: type(brown)  $\verb|nltk.corpus.reader.tagged.CategorizedTaggedCorpusReader|\\$ In [35]: **import** nltk # Download the necessary data nltk.download('abc') # Load the file from the NLTK corpus psh\_txt = nltk.corpus.abc.raw('science.txt') [nltk\_data] Downloading package abc to /home/jadi/nltk\_data... [nltk\_data] Package abc is already up-to-date! In [36]: **from** nltk.util **import** ngrams words = nltk.word\_tokenize(psh\_txt) psh\_bigrams = list(ngrams(words, 2)) psh\_trigrams = list(ngrams(words, 3)) In [39]: import nltk from nltk.corpus import brown from nltk import ngrams # Download the Brown corpus if not already downloaded nltk.download('brown') # Get the wordlist from the Brown corpus and create bigrams wordlist = brown.words() bigramlist = list(ngrams(wordlist, 2)) # Print the first 10 bigrams as an example print(bigramlist[:10]) [nltk\_data] Downloading package brown to /home/jadi/nltk\_data... [nltk\_data] Package brown is already up-to-date! [('The', 'Fulton'), ('Fulton', 'County'), ('County', 'Grand'), ('Grand', 'Jury'), ('Jury', 'said'), ('said', 'Friday', 'an'), ('an', 'investigation'), ('investigation', 'of'), ('of', "Atlanta's")] Task 3: Generating Random Text with Bigrams In [40]: cfd = nltk.ConditionalFreqDist(bigramlist) In [41]: cfd <ConditionalFreqDist with 56057 conditions> In [42]: def generate\_model(cfdist, word, num): for i in range(num): print(word, end=' ') word = cfdist[word].max() In [43]: generate\_model(cfd, 'Today', 10) Today , and the same time , and the same In [44]: generate\_model(cfd, 'Today', 30) Today , and the same time Task 30 word generated sentence In [56]: len(cfd) Out[56]: In [57]: print('30 word generated sentence') 30 word generated sentence In [58]: generate\_model(cfd, 'Today', 30) Today , and the same time , and the same