Activity • Use the harry-potter dataset Explore the pattern Identify the cleaning points Conduct text cleaning • Explore the word and character distribution Do you see different patterns from what you observed in the morning? Use the harry-potter dataset In [2]: import nltk from nltk.tokenize import word tokenize In [3]: f = open('harry 7books.txt', encoding='utf8') raw = f.read() In [4]: tokens = word tokenize(raw) In [32]: tokens[10:20] ['rowling', Out[32]: 'illustrations', 'by', 'mary', 'grandpré', 'arthur', 'a.', 'levine', 'books', 'an'l Explore the pattern In [33]: import re [w for w in tokens if re.search('ed\$', w)][10:20] ['roared', Out[33]: 'tried', 'bored', 'used', 'snarled', 'fried', 'exchanged', 'tried', 'drowned', 'drooped'] In [34]: [w for w in tokens if re.search(', |©|""|', w)][10:20] ['rowling', Out[34]: 'illustrations', 'by', 'mary', 'grandpré', 'arthur', 'a.', 'levine', 'books', 'an'] In [8]: [w for w in tokens if re.search(' $^[0-9]{1}-[0-9]{3}$ ', w)] ['0-439-06486-4', Out[8]: '0-545-02936-8', '0-439-13959-7', '0-439-78454-9', '0-439-35806-X', '0-439-13635-0', '0-590-35340-3'1 In [35]: [w for w in tokens if re.search('^[0-9]', w)][10:20] ['1999', '98-46370', '60', '59', '58', '57', '56', '55', '54', '53'] Out[35]: In [36]: [w for w in tokens if re.search('\W+\D', w)][10:20] ['mat\xadter', Out[36]: 'sleep\xading', 'four-poster', 'four\xadteen', 'top-of-the-line', 'sum\xadmer', 'wiz\xadards', 'horse-faced', 'bril\xadliant', 'jet-black'] In [37]: [w for w in tokens if re.search('\xad', w)][10:20] ['un\xadderstanding', Out[37]: 'some\xadthing', 'ig\xadnore', 'ca\xadreer', 'un\xadcle', 'talk\xading', 'wel\xadcome', 'sim\xadpering', 'compli\xadments', 'petu\xadnia'] In [38]: [w for w in tokens if re.search($'\.+$'$, w)][10:20] ['.', 'inc.', '.', 'inc.', 'bros.', 'j.', 'k.', '.', 'inc.', '.'] Out[38]: In [39]: [w for w in tokens if re.search('^[a-zA-Z]-', w)][10:20] ['c-cat', Out[39]: 'm-millicent', 'm-must', 'p-potion', 'u-bend', 'd-did', 'v-very', 'b-breakfast', 'c-couldn', 's-swear'] In [40]: [w for w in tokens if re.search('\.\$', w)][10:20] ['.', 'inc.', '.', 'inc.', 'bros.', 'j.', 'k.', '.', 'inc.', '.'] Out[40]: • Identify the cleaning points Conduct text cleaning In [41]: tokens = [w.lower() for w in tokens] tokens[10:20] ['rowling', Out[41]: 'illustrations', 'mary', 'grandpré', 'arthur', 'a.', 'levine', 'books', 'an'] In [16]: from nltk.corpus import stopwords stop words=set(stopwords.words("english")) print(stop_words) {'were', "mightn't", 'once', "aren't", 'll', 'under', 'above', "couldn't", 'ain', "it's", 'in', 'how', 'off', 're', "didn't", 'y', 'than', 'has', 'such', 'your', 'are', 'over', "wasn't", 'isn', 'should', 'that', 'didn', "shouldn't", 'why', 'can', 'his', 'mightn', 'am', 'against', 'when', 'we', 'some', 'through', 't', 'most', 'bu t', 'won', 'me', 'was', 'a', 'again', 'up', 'what', 'not', 'our', 'hadn', 'there', 'down', 'into', 'whom', 'sha n', 'him', "you've", "you'll", 'only', 'until', 'having', 'herself', 'couldn', 'between', 'haven', "needn't", 'these', "doesn't", 'be', 'needn', 'no', 'my', 'after', 'have', 'of', 'will', 'doesn', 'both', 'don', 'at', 'wh ile', 'all', 'own', 'or', 'from', 'being', 'for', 'ma', "you're", "don't", 'just', 'yours', 'nor', 'yourself', 'their', 'mustn', 'each', 'the', 'theirs', 'its', "weren't", 'she', 'other', 'few', 'he', 'here', 'i', "has n't", 'if', 've', 'ourselves', 'with', 'shouldn', 'did', 'now', 'who', 'wouldn', 'it', 'do', 'which', 'themselv es', 'her', 'very', 'too', 'those', 'then', 'them', 'o', 'during', 'further', 'as', 'on', 'any', 'hasn', 'doe s', 'to', 'an', "you'd", 'more', "wouldn't", 'they', 'm', 'hers', 'yourselves', 'where', 'been', 'because', "wo n't", "mustn't", 'you', 'itself', 'ours', 'by', 'same', "shan't", 'himself', 'before', 'and', "she's", 'so', 'b elow', 'weren', 's', 'about', 'd', "should've", "haven't", 'wasn', "that'll", 'had', 'is', 'myself', 'aren', "h adn't", 'out', "isn't", 'doing', 'this'} In [17]: filtered sent=[] for w in tokens: if w not in stop words: filtered sent.append(w) In [42]: filtered sent[10:20] ['arthur', Out[42]: 'a.', 'levine', 'books', 'imprint', 'scholastic', 'press', 'seán', 'p.', 'f.'] In [19]: from nltk.stem.wordnet import WordNetLemmatizer lem = WordNetLemmatizer() words = []for w in filtered sent: words.append(lem.lemmatize(w)) In [43]: words[10:20] ['arthur', Out[43]: 'a.', 'levine', 'book', 'imprint', 'scholastic', 'press', 'seán', 'p.', 'f.'] In [44]: clean_text = [w for w in words if re.findall('^[a-zA-Z]-', w)] +[w for w in words if re.findall('\xad', w)]+[w clean text[10:20] ['c-cat', Out[44]: 'm-millicent', 'm-must', 'p-potion', 'u-bend', 'd-did', 'v-very', 'b-breakfast', 'c-couldn', 's-swear'] In []: new_list = [w for w in words if w not in clean_text] new list[10:20] Explore the word and character distribution • Do you see different patterns from what you observed in the morning? In [28]: import matplotlib.pyplot as plt from wordcloud import WordCloud #convert list to string and generate unique_string=(" ").join(new_list) wordcloud = WordCloud(width = 1000, height = 500).generate(unique_string) In [29]: plt.imshow(wordcloud, interpolation='bilinear') plt.axis("off") plt.show() In []: