Lab1-Assignment Copyright: Vrije Universiteit Amsterdam, Faculty of Humanities, CLTL This notebook describes the assignment for Lab 1 of the text mining course. Points: each exercise is prefixed with the number of points you can obtain for the exercise. We assume you have worked through the following notebooks: Lab1.1-introduction Lab1.2-introduction-to-NLTK Lab1.3-introduction-to-spaCy In this assignment, you will process an English text (Lab1-apple-samsung-example.txt) with both NLTK and spaCy and discuss the similarities and differences. **Credits** The notebooks in this block have been originally created by Marten Postma. Adaptations were made by Filip Ilievski. Tip: how to read a file from disk Let's open the file Lab1-apple-samsung-example.txt from disk. In [3]: from pathlib import Path In [4]: | nltk.download('averaged\_perceptron\_tagger') nltk.download('maxent ne chunker') Traceback (most recent call last) NameError <ipython-input-4-b3ef775cd85f> in <module> ---> 1 nltk.download('averaged\_perceptron\_tagger') 2 nltk.download('maxent ne chunker') NameError: name 'nltk' is not defined In [5]: cur\_dir = Path().resolve() # this should provide you with the folder in which this not path\_to\_file = Path.joinpath(cur\_dir, 'Lab1-apple-samsung-example.txt') print(path\_to\_file) print('does path exist? ->', Path.exists(path\_to\_file)) C:\Users\Marlon\Downloads\Lab1-apple-samsung-example.txt does path exist? -> True If the output from the code cell above states that does path exist? -> False, please check that the file **Lab1-apple-samsung-example.txt** is in the same directory as this notebook. with open (path to file) as infile: text = infile.read() print('number of characters', len(text)) number of characters 1142 [total points: 4] Exercise 1: NLTK In this exercise, we use NLTK to apply Part-of-speech (POS) tagging, Named Entity Recognition (NER), and Constituency parsing. The following code snippet already performs sentence splitting and tokenization. import nltk from nltk.tokenize import sent tokenize from nltk import word\_tokenize sentences\_nltk = sent\_tokenize(text) In [8]: tokens per sentence = [] for sentence nltk in sentences nltk: sent tokens = word tokenize(sentence nltk) tokens per sentence.append(sent tokens) We will use lists to keep track of the output of the NLP tasks. We can hence inspect the output for each task using the index of the sentence.  $sent_id = 1$ print('SENTENCE', sentences\_nltk[sent\_id]) print('TOKENS', tokens\_per\_sentence[sent id]) SENTENCE The six phones and tablets affected are the Galaxy S III, running the new Jel ly Bean system, the Galaxy Tab 8.9 Wifi tablet, the Galaxy Tab 2 10.1, Galaxy Rugby Pr o and Galaxy S III mini. TOKENS ['The', 'six', 'phones', 'and', 'tablets', 'affected', 'are', 'the', 'Galaxy', 'S', 'III', ',', 'running', 'the', 'new', 'Jelly', 'Bean', 'system', ',', 'the', 'Galaxy', 'Tab', '8.9', 'Wifi', 'tablet', ',', 'the', 'Galaxy', 'Tab', '2', '10.1', ',', 'Galaxy', 'Rugby', 'Pro', 'and', 'Galaxy', 'S', 'III', 'mini', '.'] [point: 1] Exercise 1a: Part-of-speech (POS) tagging Use nltk.pos\_tag to perform part-of-speech tagging on each sentence. Use print to **show** the output in the notebook (and hence also in the exported PDF!). In [11]: pos tags per sentence = [] for tokens in tokens per sentence: print(nltk.pos tag(tokens)) pos tags per sentence.append(nltk.pos tag(tokens)) [('https', 'NN'), (':', ':'), ('//www.telegraph.co.uk/technology/apple/9702716/Apple-S amsung-lawsuit-six-more-products-under-scrutiny.html', 'JJ'), ('Documents', 'NNS'), ('filed', 'VBN'), ('to', 'TO'), ('the', 'DT'), ('San', 'NNP'), ('Jose', 'NNP'), ('fede ral', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Nov ember', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('products', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('`', '`'), ('Jelly', 'RB'), ('Be an', 'NNP'), ("''", "''"), ('and', 'CC'), ('`'', '`''), ('Ice', 'NNP'), ('Cream', 'NN P'), ('Sandwich', 'NNP'), ("''", "''"), ('Operating', 'VBG'), ('systems', 'NNS') P'), ('Sandwich', 'NNP'), ("''", "''"), ('operating', 'VBG'), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), ('claims', 'VBZ'), ('infringe', 'VB'), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')]
[('The', 'DT'), ('six', 'CD'), ('phones', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), ('affected', 'VBN'), ('are', 'VBP'), ('the', 'DT'), ('Galaxy', 'NNP'), ('III', 'NNP'), (',',','), ('running', 'VBG'), ('the', 'DT'), ('new', 'JJ'), ('Jell y', 'NNP'), ('Bean', 'NNP'), ('system', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('8.9', 'CD'), ('Wifi', 'NNP'), ('tablet', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), ('Galaxy', 'NNP'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), ('mini', 'NN'), ('.', '.')]
[('Apple', 'NNP'), ('stated', 'VBD'), ('it', 'PRP'), ('had', 'VBD'), ('"acted', 'VBN'), ('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ("''", "''"), ('in', 'IN'), ('order', 'NN'), ('to', 'TO'), ('``', '``'), ('determine', 'VB'), ('that', 'IN'), ('these', 'DT'), ('newly', 'RB'), ('released', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'JJ'), ('of', 'IN'), ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('already', 'RB'), ('asserted', 'VBN'), ('by', 'IN'), ('Apple', 'NN P'), ('.', '.'), ("''", "''")] P'), ('.', '.'), ("''", "''")]
[('In', 'IN'), ('August', 'NNP'), (',', ','), ('Samsung', 'NNP'), ('lost', 'VBD'),
('a', 'DT'), ('US', 'NNP'), ('patent', 'NN'), ('case', 'NN'), ('to', 'TO'), ('Apple',
'NNP'), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('to', 'TO'), ('pay', 'V

B'), ('its', 'PRP\$'), ('rival', 'JJ'), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), ('Â

£0.66bn', 'NN'), (')', ')'), ('in', 'IN'), ('damages', 'NNS'), ('for', 'IN'), ('copyin
g', 'VBG'), ('features', 'NNS'), ('of', 'IN'), ('the', 'DT'), ('iPad', 'NN'), ('and',
'CC'), ('iPhone', 'NN'), ('in', 'IN'), ('its', 'PRP\$'), ('Galaxy', 'NNP'), ('range',
'NN'), ('of', 'IN'), ('devices', 'NNS'), ('.', '.')]
[('Samsung', 'NNP'), (',', ','), ('which', 'WDT'), ('is', 'VBZ'), ('the', 'DT'), ('wor
ld', 'NN'), ("'s", 'POS'), ('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker',
'NN'), (',', ','), ('is', 'VBZ'), ('appealing', 'VBG'), ('the', 'DT'), ('ruling', 'N
N'), ('.', '.')]
[('A', 'DT'), ('similar', 'JJ'), ('case', 'NN'), ('in', 'IN'), ('the', 'DT'), ('UK', [('A', 'DT'), ('similar', 'JJ'), ('case', 'NN'), ('in', 'IN'), ('the', 'DT'), ('UK', 'NNP'), ('found', 'VBD'), ('in', 'IN'), ('Samsung', 'NNP'), ("'s", 'POS'), ('favour', 'NNP'), ('and', 'CC'), ('ordered', 'VBD'), ('Apple', 'NNP'), ('to', 'TO'), ('publish', 'VB'), ('an', 'DT'), ('apology', 'NN'), ('making', 'VBG'), ('clear', 'JJ'), ('that', 'IN'), ('the', 'DT'), ('South', 'JJ'), ('Korean', 'JJ'), ('firm', 'NN'), ('had', 'VBD'), ('not', 'RB'), ('copied', 'VBN'), ('its', 'PRP\$'), ('iPad', 'NN'), ('when', 'WRD'), ('designing', 'VBG'), ('its', 'PRP\$'), ('own', 'JJ'), ('devices', 'NNS'), ('.', 'NNS'), ( In [12]: print(pos tags per sentence) [[('https', 'NN'), (':', ':'), ('//www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more-products-under-scrutiny.html', 'JJ'), ('Documents', 'NNS'), ('filed', 'VBN'), ('to', 'TO'), ('the', 'DT'), ('San', 'NNP'), ('Jose', 'NNP'), ('fede ral', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Nov ember', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('products', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('`', '`'), ('Jelly', 'RB'), ('Be an', 'NNP'), ("''", "''"), ('and', 'CC'), ('`'', '`''), ('Ice', 'NNP'), ('Cream', 'NN P'), ('Sandwich', 'NNP'), ("''", "''"), ('soperating', 'VBG'), ('systems', 'NNS') an', 'NNP'), ("'", "'"), ('and', 'CC'), ('``', '``'), ('Ice', 'NNP'), ('Cream', 'NN P'), ('Sandwich', 'NNP'), ("'", "'"), ('operating', 'VBG'), ('systems', 'NNS'), (',',','), ('which', 'WDT'), ('Apple', 'NNP'), ('claims', 'VBZ'), ('infringe', 'VB'), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')], [('The', 'DT'), ('six', 'CD'), ('pho nes', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), ('affected', 'VBN'), ('are', 'VBP'), ('the', 'DT'), ('Galaxy', 'NNP'), ('III', 'NNP'), (',', ','), ('runnin g', 'VBG'), ('the', 'DT'), ('new', 'JJ'), ('Jelly', 'NNP'), ('Bean', 'NNP'), ('syste m', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('8.9', 'C D'), ('Wifi', 'NNP'), ('tablet', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), ('Galaxy', 'NNP'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NN P'), ('mini', 'NN'), ('.', '.')], [('Apple', 'NNP'), ('stated', 'VBD'), ('it', 'PRP'), ('had', 'VBD'), ('â@acted', 'VBN'), ('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ("""", """"), ('in', 'IN'), ('order', 'NN'), ('to', 'TO'), ('`', '`'), ('dete rmine', 'VB'), ('that', 'IN'), ('these', 'DT'), ('newly', 'RB'), ('released', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'JJ'), ('of', 'IN'), ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('already', 'RB'), ('asserted', 'VB ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('already', 'RB'), ('asserted', 'VB N'), ('by', 'IN'), ('Apple', 'NNP'), ('.', '.'), ("''", "''")], [('In', 'IN'), ('Augus t', 'NNP'), (',', ','), ('Samsung', 'NNP'), ('lost', 'VBD'), ('a', 'DT'), ('US', 'NN P'), ('patent', 'NN'), ('case', 'NN'), ('to', 'TO'), ('Apple', 'NNP'), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('to', 'TO'), ('pay', 'VB'), ('its', 'PRP\$'), ('ri val', 'JJ'), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), ('£0.66bn', 'NN'), (')', ')'), val', 'JJ'), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), ('A£0.66bn', 'NN'), (')', ')'),
('in', 'IN'), ('damages', 'NNS'), ('for', 'IN'), ('copying', 'VBG'), ('features', 'NN
S'), ('of', 'IN'), ('the', 'DT'), ('iPad', 'NN'), ('and', 'CC'), ('iPhone', 'NN'), ('i
n', 'IN'), ('its', 'PRP\$'), ('Galaxy', 'NNP'), ('range', 'NN'), ('of', 'IN'), ('device
s', 'NNS'), ('.', '.')], [('Samsung', 'NNP'), (',', ','), ('which', 'WDT'), ('is', 'VB
Z'), ('the', 'DT'), ('world', 'NN'), ("'s", 'POS'), ('top', 'JJ'), ('mobile', 'NN'),
('phone', 'NN'), ('maker', 'NN'), (',', ','), ('is', 'VBZ'), ('appealing', 'VBG'), ('t
he', 'DT'), ('ruling', 'NN'), ('.', '.')], [('A', 'DT'), ('similar', 'JJ'), ('case',
'NN'), ('in', 'IN'), ('the', 'DT'), ('UK', 'NNP'), ('found', 'VBD'), ('in', 'IN'), ('s
amsung', 'NNP'), ("'s", 'POS'), ('favour', 'NN'), ('and', 'CC'), ('ordered', 'VBD') amsung', 'NNP'), ("'s", 'POS'), ('favour', 'NN'), ('and', 'CC'), ('ordered', 'VBD'), ('Apple', 'NNP'), ('s, 'ros'), ('favour', 'NN'), ('and', 'CC'), ('ordered', 'VBD'), ('Apple', 'NNP'), ('to', 'TO'), ('publish', 'VB'), ('an', 'DT'), ('apology', 'NN'), ('making', 'VBG'), ('clear', 'JJ'), ('that', 'IN'), ('the', 'DT'), ('South', 'JJ'), ('Korean', 'JJ'), ('firm', 'NN'), ('had', 'VBD'), ('not', 'RB'), ('copied', 'VBN'), ('its', 'PRP\$'), ('iPad', 'NN'), ('when', 'WRB'), ('designing', 'VBG'), ('its', 'PRP\$'), ('own', 'JJ'), ('devices', 'NNS'), ('.', '.')]] [point: 1] Exercise 1b: Named Entity Recognition (NER) Use nltk.chunk.ne\_chunk to perform Named Entity Recognition (NER) on each sentence. Use print to **show** the output in the notebook (and hence also in the exported PDF!). In [13]: **from** nltk.chunk **import** ne chunk ner\_tags\_per\_sentence = [] for sentence in pos\_tags\_per\_sentence: ner tags per sentence.append(nltk.chunk.ne chunk(sentence)) In [14]: print(ner tags per sentence) [Tree('S', [('https', 'NN'), (':', ':'), ('//www.telegraph.co.uk/technology/apple/9702 716/Apple-Samsung-lawsuit-six-more-products-under-scrutiny.html', 'JJ'), ('Documents', 'NNS'), ('filed', 'VBN'), ('to', 'TO'), ('the', 'DT'), Tree('ORGANIZATION', [('San', 'NNP'), ('Jose', 'NNP')]), ('federal', 'JJ'), ('court', 'NN'), ('in', 'IN'), Tree('GP E', [('California', 'NNP')]), ('on', 'IN'), ('November', 'NNP'), ('23', 'CD'), ('lis t', 'NN'), ('six', 'CD'), Tree('ORGANIZATION', [('Samsung', 'NNP')]), ('products', 'NN S'), ('running', 'VBG'), ('the', 'DT'), ('``', '``'), ('Jelly', 'RB'), Tree('GPE', [('Bean', 'NNP')]), ("''", "''"), ('and', 'CC'), ('``', '``'), ('Ice', 'NNP'), ('Crea [('Bean', 'NNP')]), ('and, 'cc'), ('ind', 'cc'), ('ind', 'cc'), ('cfea', 'NNP')], ('systems', 'NN 'nn'), ('systems', 'NN'), (',', ','), ('which', 'WDT'), Tree('PERSON', [('Apple', 'NNP')]), ('claims', 'VB 'z'), ('infringe', 'VB'), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')]), Tree('S', [('The', 'DT'), ('six', 'CD'), ('phones', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), ('affected', 'VBN'), ('are', 'VBP'), ('the', 'DT'), Tree('ORGANIZATION', [('Galaxy', 'NNP')]), ('S', 'NNP'), ('III', 'NNP'), (',', ','), ('running', 'VBG'), ('the', 'DT'), ('pay', 'III')  $\hbox{('new', 'JJ'), Tree('PERSON', [('Jelly', 'NNP'), ('Bean', 'NNP')]), ('system', 'NN'), }$ (',',','), ('the', 'DT'), Tree('ORGANIZATION', [('Galaxy', 'NNP')]), ('Tab', 'NNP'), ('8.9', 'CD'), ('Wifi', 'NNP'), ('tablet', 'NN'), (',', ','), ('the', 'DT'), Tree('ORG ANIZATION', [('Galaxy', 'NNP')]), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), Tree('PERSON', [('Galaxy', 'NNP'), ('Rugby', 'NNP'), ('Pro', 'NNP')]), ('and', 'CC'), Tree('PERSON', [('Galaxy', 'NNP'), ('S', 'NNP')]), ('III', 'NNP'), ('mini', 'NN'), ('.', '.')]), Tree('S', [Tree('PERSON', [('Apple', 'NNP')]), ('stated', 'VBD'), ('it', 'PRP'), ('had', 'VBD'), ('"acted', 'VBN'), ('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ("'", ""'"), ('in', 'IN'), ('order', 'NN'), ('to', 'TO'), ('`', '``'), ('determine', 'VB'), ('that', 'IN'), ('these', 'DT'), ('newly', 'RB'), ('reelessed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'Interpressed', 'VBN'), ('ma leased', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'J J'), ('of', 'IN'), ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('already', 'R B'), ('asserted', 'VBN'), ('by', 'IN'), Tree('PERSON', [('Apple', 'NNP')]), ('.', '.'), ("''", "''")]), Tree('S', [('In', 'IN'), Tree('GPE', [('August', 'NNP')]), (',', ','), Tree('PERSON', [('Samsung', 'NNP')]), ('lost', 'VBD'), ('a', 'DT'), Tree('GSP', [('US', 'NNP')]), ('patent', 'NN'), ('case', 'NN'), ('to', 'TO'), Tree('GPE', [('Apple', 'NNP')]), ('and 'prod', 'VBN'), ('and 'prod', 'VBN'), ('and 'prod', 'NNP')]), ('patent', 'NNP'), ('prodered', 'VBN'), ('to', 'TO'), Tree('TPR'), ('prodered', 'VBN'), ('to', 'TO'), 'Tree', 'TO'), ('prod', 'VBN'), ('prod', 'VB e', 'NNP')]), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('to', 'TO'), ('pay', 'VB'), ('its', 'PRP\$'), ('rival', 'JJ'), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), ('Â'))) £0.66bn', 'NN'), (')', ')'), ('in', 'IN'), ('damages', 'NNS'), ('for', 'IN'), ('copyin g', 'VBG'), ('features', 'NNS'), ('of', 'IN'), ('the', 'DT'), Tree('ORGANIZATION', ('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker', 'NN'), (',', ','), ('is',
'VBZ'), ('appealing', 'VBG'), ('the', 'DT'), ('ruling', 'NN'), ('.', '.')]), Tree('S',
[('A', 'DT'), ('similar', 'JJ'), ('case', 'NN'), ('in', 'IN'), ('the', 'DT'), Tree('OR
GANIZATION', [('UK', 'NNP')]), ('found', 'VBD'), ('in', 'IN'), Tree('GPE', [('Samsun
g', 'NNP')]), ("'s", 'POS'), ('favour', 'NN'), ('and', 'CC'), ('ordered', 'VBD'), Tree
('PERSON', [('Apple', 'NNP')]), ('to', 'TO'), ('publish', 'VB'), ('an', 'DT'), ('apolo
gy', 'NN'), ('making', 'VBG'), ('clear', 'JJ'), ('that', 'IN'), ('the', 'DT'), Tree('L
OCATION', [('South', 'JJ'), ('Korean', 'JJ')]), ('firm', 'NN'), ('had', 'VBD'), ('no
t', 'RB'), ('copied', 'VBN'), ('its', 'PRP\$'), ('ipad', 'NN'), ('when', 'WRB'), ('desi
gning', 'VBG'), ('its', 'PRP\$'), ('devices', 'NNS'), (''', '')])] gning', 'VBG'), ('its', 'PRP\$'), ('own', 'JJ'), ('devices', 'NNS'), ('.', '.')])] [points: 2] Exercise 1c: Constituency parsing Use the nltk.RegexpParser to perform constituency parsing on each sentence. Use print to **show** the output in the notebook (and hence also in the exported PDF!). In [16]: constituent parser = nltk.RegexpParser(''' NP:  ${<DT>? < JJ>* < NN>*}$  # NP P: {<IN>} # Preposition V: {<V.\*>} PP: {<P> <NP>} # PP -> P NP VP: {<V> <NP|PP>\*} # VP -> V (NP|PP)\*''') constituency\_output\_per\_sentence = [] for token in pos\_tags\_per\_sentence: constituency\_output\_per\_sentence.append(constituent\_parser.parse(token)) print(constituency\_output\_per\_sentence) o.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more-products-under-scrutiny.h tml', 'JJ')]), ('Documents', 'NNS'), Tree('VP', [Tree('V', [('filed', 'VBN')])]), ('t o', 'TO'), Tree('NP', [('the', 'DT')]), ('San', 'NNP'), ('Jose', 'NNP'), Tree('NP', [('federal', 'JJ'), ('court', 'NN')]), Tree('P', [('in', 'IN')]), ('California', 'NN 'VBG')])]), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), Tree ('VP', [Tree('V', [('claims', 'VBZ')])]), Tree('VP', [Tree('V', [('infringe', 'VB')])]), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')]), Tree('S', [Tree('NP', [('T he', 'DT')]), ('six', 'CD'), ('phones', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), Tre e('VP', [Tree('V', [('affected', 'VBN')])]), Tree('VP', [Tree('V', [('are', 'VBP')]), Tree('NP', [('the', 'DT')])]), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (',', ','), Tree('VP', [Tree('V', [('running', 'VBG')]), Tree('NP', [('the', 'DT'), ('new', 'III')]))), ('Tree('VP', [Tree('V', [('running', 'VBG')]), Tree('NP', [('the', 'DT'), ('new', 'III')]))) 'JJ')])), ('Jelly', 'NNP'), ('Bean', 'NNP'), Tree('NP', [('the', 'DT'), ('new', 'JJ')]))), ('Jelly', 'NNP'), ('Bean', 'NNP'), Tree('NP', [('system', 'NN')]), (',', ','), Tree('NP', [('the', 'DT')]), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('Tab', 'NNP'), (',', ','), Tree('NP', [('the', 'DT')]), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), ('Galaxy', 'NNP'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), Tree('NP', [('mini', 'NN')]), ('.', '.')]), Tree('S', [('Apple', 'NNP'), Tree('VP', [Tree('V', [('stated', 'VBD')])]), ('it', 'PRP'), Tree [('Apple', 'NNP'), Tree('VP', [Tree('V', [('stated', 'VBD')])]), ('it', 'PRP'), Tree
('VP', [Tree('V', [('had', 'VBD')])]), Tree('VP', [Tree('V', [('â@acted', 'VBN')])]),
('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ("''", "''"), Tree('PP', [Tree
('P', [('in', 'IN')]), Tree('NP', [('order', 'NN')])]), ('to', 'TO'), ('``', '``'), Tr
ee('VP', [Tree('V', [('determine', 'VB')]), Tree('PP', [Tree('P', [('that', 'IN')]), Tree('NP', [('these', 'DT')])])]), ('newly', 'RB'), Tree('VP', [Tree('V', [('do', 'VBP')])]), Tree('VP', [Tree('V', [('infringe', 'VB')]), Tree('NP', [('many', 'JJ')]), Tree('PP', [Tree('P', [('of', 'IN')]), Tree('NP', [('the', 'DT'), ('same', 'JJ')])])), Tree('PP', [('by', 'laready', 'RB'), Tree('VP', [Tree('V', [('asserted', 'VBN')])]), Tree('PP', [('by', 'laready', 'RB'), Tree('VP', [('by', 'VBN')])]), Tree('PP', [('by', 'VBN')])]) ('already', 'RB'), Tree('VP', [Tree('V', [('asserted', 'VBN')])]), Tree('P', [('by', Tree('P', Tree('P', Tree('N', Tree('N','IN')]), ('Apple', 'NNP'), ('.', '.'), ("''", "''")]), Tree('S', [Tree('P', [('In', 'I N')]), ('August', 'NNP'), (',', ','), ('Samsung', 'NNP'), Tree('VP', [Tree('V', [('los t', 'VBD')]), Tree('NP', [('a', 'DT')])]), ('US', 'NNP'), Tree('NP', [('patent', 'N N'), ('case', 'NN')]), ('to', 'TO'), ('Apple', 'NNP'), ('and', 'CC'), Tree('VP', [Tree ('V', [('was', 'VBD')])]), Tree('VP', [Tree('V', [('ordered', 'VBN')])]), ('to', 'T O'), Tree('VP', [Tree('V', [('pay', 'VB')])]), ('its', 'PRP\$'), Tree('NP', [('rival', 'JJ')]), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), Tree('NP', [('£0.66bn', 'NN')]), (')', ')'), Tree('P', [('in', 'IN')]), ('damages', 'NNS'), Tree('P', [('for', 'IN')]), Tree('VP', [Tree('V', [('copying', 'VBG')])]), ('features', 'NNS'), Tree('PP', [Tree ('P', [('of', 'IN')]), Tree('NP', [('the', 'DT'), ('iPad', 'NN')])]), ('and', 'CC'), Tree('NP', [('iPhone', 'NN')]), Tree('P', [('in', 'IN')]), ('its', 'PRP\$'), ('Galaxy', 'NNP'), Tree('NP', [('range', 'NN')]), Tree('P', [('of', 'IN')]), ('devices', 'NNS'), ('.', '.')]), Tree('S', [('Samsung', 'NNP'), (',', ','), ('which', 'WDT'), Tree('VP', [Tree('V', [('is', 'VBZ')]), Tree('NP', [('the', 'DT'), ('world', 'NN')])]), ("'s", 'PROS'), Tree('NP', [('the', 'DT'), ('world', 'NN')])]), Tree('NP', [('the', 'DT'), ('world', 'NN')])]), ("'s", 'PROS'), Tree('NP', [('the', 'DT'), ('world', 'NN')])]), Tree('NN'), ('world', 'NN')])] OS'), Tree('NP', [('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker', 'NN')]), (',', ','), Tree('VP', [Tree('V', [('is', 'VBZ')])]), Tree('VP', [Tree('V', [('appealing', 'VBG')]), Tree('NP', [('the', 'DT'), ('ruling', 'NN')])]), ('.', '.')]), Tree ('S', [Tree('NP', [('A', 'DT'), ('similar', 'JJ'), ('case', 'NN')]), Tree('PP', [Tree ('P', [('in', 'IN')]), Tree('NP', [('the', 'DT')])]), ('UK', 'NNP'), Tree('VP', [Tree ('V', [('found', 'VBD')])]), Tree('P', [('in', 'IN')]), ('Samsung', 'NNP'), ("'s", 'POS'), Tree('NP', [('favour', 'NN')]), ('and', 'CC'), Tree('VP', [Tree('V', [('ordered', 'VBD')])]), ('Apple', 'NNP'), ('to', 'TO'), Tree('VP', [Tree('V', [('publish', 'VBD')]), Tree('NP', [('an', 'DT'), ('apology', 'NN')])]), Tree('VP', [Tree('V', [('that', 'IN')]), Tree('NP', [('the', 'DT'), ('South', 'JJ'), ('Korean', 'JJ'), ('firm', 'NN')])])]), Tree('VP', [Tree('V', [('that', 'IN')])]), Tree('VP', [Tree('V', [('that', 'IN')])]), Tree('VP', [Tree('V', [('that', 'NN')])])]), Tree('VP', [Tree('V', [('that', 'NN')])])]), Tree('VP', [Tree('V', [('that', 'NN')])])]), Tree('VP', [Tree('V', [('that', 'NN')])])]), Tree('VP', [Tree('V', [('topie', 'NN')])]), ('not', 'RB'), Tree('VP', [Tree('V', [('copie', 'NN')])]), ('not', 'RB'), Tree('VP', [Tree('V', [('copie', 'NN')])])]), ('not', 'RB'), Tree('VP', [Tree('V', [('copie', 'NN')])])]), ('not', 'RB'), Tree('VP', [Tree('V', [('copie', 'NN')])])]) OS'), Tree('NP', [('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker', 'NN')]), ee('VP', [Tree('V', [('had', 'VBD')])]), ('not', 'RB'), Tree('VP', [Tree('V', [('copie d', 'VBN')]))), ('its', 'PRP\$'), Tree('NP', [('iPad', 'NN')]), ('when', 'WRB'), Tree
('VP', [Tree('V', [('designing', 'VBG')]))), ('its', 'PRP\$'), Tree('NP', [('own', 'J J')]), ('devices', 'NNS'), ('.', '.')])] Augment the RegexpParser so that it also detects Named Entity Phrases (NEP), e.g., that it detects Galaxy S III and Ice Cream Sandwich constituent\_parser\_v2 = nltk.RegexpParser('''  $NP: \{ < DT > ? < JJ > * < NN > * \} # NP$ P: {<IN>} # Preposition V: {<V.\*>} # Verb # PP -> P NP PP: {<P> <NP>} NEP: {<NNP>\*} # ???''') constituency\_v2\_output\_per\_sentence = [] for token in pos\_tags\_per\_sentence: constituency v2 output per sentence.append(constituent parser.parse(token)) print(constituency\_v2\_output\_per\_sentence) o.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more-products-under-scrutiny.h tml', 'JJ')]), ('Documents', 'NNS'), Tree('VP', [Tree('V', [('filed', 'VBN')])]), ('t
o', 'TO'), Tree('NP', [('the', 'DT')]), ('San', 'NNP'), ('Jose', 'NNP'), Tree('NP',
[('federal', 'JJ'), ('court', 'NN')]), Tree('P', [('in', 'IN')]), ('California', 'NN P'), Tree('P', [('on', 'IN')]), ('November', 'NNP'), ('23', 'CD'), Tree('NP', [('lis t', 'NN')]), ('six', 'CD'), ('Samsung', 'NNP'), ('products', 'NNS'), Tree('VP', [Tree t', 'NN')]) c , NNN')]), ('Six', 'CD'), ('Samsung', 'NNP'), ('products', 'NNS'), Tree('VP', [Tree
('V', [('running', 'VBG')]), Tree('NP', [('the', 'DT')])]), ('``', '``'), ('Jelly', 'R
B'), ('Bean', 'NNP'), ("''", "''"), ('and', 'CC'), ('``', '``'), ('Ice', 'NNP'), ('Cre
am', 'NNP'), ('Sandwich', 'NNP'), ("''", "''"), Tree('VP', [Tree('V', [('operating',
'VBG')])]), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), Tree
('VP', [Tree('V', [('claims', 'VBZ')])]), Tree('VP', [Tree('V', [('infringe', 'V
B')])]), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')]), Tree('S', [Tree('NP', [('Tee('VP', [Tree('V', [('are', 'VRP')])]), Tree('VP', [Tree('V', [('are', 'VRP')])]) e('VP', [Tree('V', [('affected', 'VBN')])]), Tree('VP', [Tree('V', [('are', 'VBP')]), Tree('NP', [('the', 'DT')])]), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (',', ','), Tree('VP', [Tree('V', [('running', 'VBG')]), Tree('NP', [('the', 'DT'), ('new', ','), Tree('VP', [Tree('V', [('running', 'VBG')]), Ifee('NF', [('che', 'DI'), ('new', 'JJ')])]), ('Jelly', 'NNP'), ('Bean', 'NNP'), Tree('NP', [('system', 'NN')]), (',', ','), Tree('NP', [('the', 'DT')]), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('Tab', 'NNP'), (',', ','), Tree('NP', [('the', 'DT')]), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), ('Galaxy', 'NNP'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('C', 'NNP'), ('C' ('S', 'NNP'), ('III', 'NNP'), Tree('NP', [('mini', 'NN')]), ('.', '.')]), Tree('S', [('Apple', 'NNP'), Tree('VP', [Tree('V', [('stated', 'VBD')])]), ('it', 'PRP'), Tree [('Apple', 'NNP'), Tree('VP', [Tree('V', [('stated', 'VBD')])]), ('it', 'PRP'), Tree
('VP', [Tree('V', [('had', 'VBD')])]), Tree('VP', [Tree('V', [('â@acted', 'VBN')])]),
('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ("''", "''"), Tree('PP', [Tree
('P', [('in', 'IN')]), Tree('NP', [('order', 'NN')])]), ('to', 'TO'), ('``', '``'), Tr
ee('VP', [Tree('V', [('determine', 'VB')]), Tree('PP', [Tree('P', [('that', 'IN')]), Tree('NP', [('these', 'DT')])])]), ('newly', 'RB'), Tree('VP', [Tree('V', [('released', 'VBN')])]), ('products', 'NNS'), Tree('VP', [('many', 'JJ')]), Tree('PP', [Tree('P', [('of', 'IN')]), Tree('NP', [('the', 'DT'), ('same', 'JJ')])])), ('claims', 'NNS'),
('claready', 'PR'), Tree('VP', [Tree('V', [('asserted', 'VBN')])]), Tree('PV', [('by', 'Intered', 'VBN', 'INNS')])), Tree('PV', [('by', 'Intered', 'VBN', 'INNS')])]), Tree('PV', [('by', 'Intered', 'VBN', 'INNS')]), Tree('PV', [('by', 'Intered', 'VBN', 'INNS')])]), Tree('PV', [('by', 'Intered', 'VBN', 'INNS', 'Intered', 'VBN', 'INNS')])]), Tree('PV', [('by', 'Intered', 'VBN', 'Intered' [('or', 'ln')]), Tree('NP', [('the', 'DT'), ('same', 'JJ')])])]), ('claims', 'NNS'),
('already', 'RB'), Tree('VP', [Tree('V', [('asserted', 'VBN')])]), Tree('P', [('by',
'IN')]), ('Apple', 'NNP'), ('.', '.'), ("''", "''")]), Tree('S', [Tree('P', [('In', 'I
N')]), ('August', 'NNP'), (',', ','), ('Samsung', 'NNP'), Tree('VP', [Tree('V', [('los
t', 'VBD')]), Tree('NP', [('a', 'DT')])]), ('US', 'NNP'), Tree('NP', [('patent', 'N
N'), ('case', 'NN')]), ('to', 'TO'), ('Apple', 'NNP'), ('and', 'CC'), Tree('VP', [Tree
('V', [('was', 'VBD')])]), Tree('VP', [Tree('V', [('ordered', 'VBN')])]), ('to', 'T
O'), Tree('VP', [Tree('V', [('pay', 'VB')])]), ('its', 'PRP\$'), Tree('NP', [('rival',
'JJ')]), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), Tree('NP', [('for', 'IN')]),
(')', ')'), Tree('P', [('in', 'IN')]), ('damages', 'NNS'), Tree('P', [('for', 'IN')]), (')', ')'), Tree('P', [('in', 'IN')]), ('damages', 'NNS'), Tree('P', [('for', 'IN')]), Tree('VP', [Tree('V', [('copying', 'VBG')])]), ('features', 'NNS'), Tree('PP', [Tree ('VP', [('of', 'IN')]), Tree('NP', [('the', 'DT'), ('iPad', 'NN')])]), ('and', 'CC'), Tree('NP', [('iPhone', 'NN')]), Tree('P', [('in', 'IN')]), ('its', 'PRP\$'), ('Galaxy', 'NNP'), Tree('NP', [('range', 'NN')]), Tree('P', [('of', 'IN')]), ('devices', 'NNS'), ('.', '.')]), Tree('S', [('Samsung', 'NNP'), (',', ','), ('which', 'WDT'), Tree('VP', [Tree('V', [('is', 'VBZ')]), Tree('NP', [('the', 'DT'), ('world', 'NN')])]), ("'s", 'P OS'), Tree('NP', [('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker', 'NN')]), OS'), Tree('NP', [('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker', 'NN')]), (',', ','), Tree('VP', [Tree('V', [('is', 'VBZ')])]), Tree('VP', [Tree('V', [('appealing', 'VBG')]), Tree('NP', [('the', 'DT'), ('ruling', 'NN')])]), ('.', '.')]), Tree ('S', [Tree('NP', [('A', 'DT'), ('similar', 'JJ'), ('case', 'NN')]), Tree('PP', [Tree ('P', [('in', 'IN')]), Tree('NP', [Tree ('V', [('found', 'VBD')])]), Tree('NP', [('in', 'IN')]), ('Samsung', 'NNP'), ("'s", 'PO S'), Tree('NP', [('favour', 'NN')]), ('and', 'CC'), Tree('VP', [Tree('V', [('ordered', 'VBD')])]), ('Apple', 'NNP'), ('to', 'TO'), Tree('VP', [Tree('V', [('publish', 'V B')])]), Tree('NP', ['an', 'DT'), ('analogy', 'NN')])]), Tree('NP', [Tree('V', [('making', 'NN')]), Tree('V', [(' B')]), Tree('NP', [('an', 'DT'), ('apology', 'NN')])]), Tree('VP', [Tree('V', [('makin g', 'VBG')]), Tree('NP', [('clear', 'JJ')]), Tree('PP', [Tree('P', [('that', 'IN')]), Tree('NP', [('the', 'DT'), ('South', 'JJ'), ('Korean', 'JJ'), ('firm', 'NN')])])]), Tr ee('VP', [Tree('V', [('had', 'VBD')])]), ('not', 'RB'), Tree('VP', [Tree('V', [('copie d', 'VBN')])), ('its', 'PRP\$'), Tree('NP', [('iPad', 'NN')]), ('when', 'WRB'), Tree
('VP', [Tree('V', [('designing', 'VBG')])]), ('its', 'PRP\$'), Tree('NP', [('own', 'J J')]), ('devices', 'NNS'), ('.', '.')])] [total points: 1] Exercise 2: spaCy Use Spacy to process the same text as you analyzed with NLTK. In [21]: import spacy nlp = spacy.load('en core web sm') doc = nlp(text) # insert code here pos\_tags\_per\_sentence\_spacy = [] ner\_tags\_per\_sentence\_spacy = [] constituency\_output\_per\_sentence\_spacy = [] for sentence in doc.sents: POS = [(token.text,token.tag\_) for token in sentence] NER = [ (ent.text,ent.label\_) for ent in sentence.ents] constituency = [(token.text, token.dep\_, token.head) for token in sentence] pos\_tags\_per\_sentence\_spacy.append(POS) ner\_tags\_per\_sentence\_spacy.append(NER) constituency\_output\_per\_sentence\_spacy.append(constituency) print(pos tags per sentence spacy) [[('https://www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-mor e-products-under-scrutiny.html', 'NNP'), (' $\n'$ , '\_SP'), ('Documents', 'NNS'), ('file 'DT'), ('San', 'NNP'), ('Jose', 'VBD'), ('to' 'IN'), ('the', 'NNP'), ('federal', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Novembe r', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('product s', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('"', '`'), ('Jelly', 'NNP'), ('Bean', 'NNP'), ('"', "''"), ('and', 'CC'), ('"', '``'), ('Ice', 'NNP'), ('Cream', 'NNP'), ('S andwich', 'NNP'), ('"', "''"), ('operating', 'NN'), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), ('claims', 'NNS'), ('infringe', 'VBP'), ('its', RP\$'), ('patents', 'NNS'), ('.', '.')], [('\n', '\_SP')], [('The', 'DT'), ('six', 'CD'), ('phones', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), ('affected', 'VBN'), ('are', 'VBP'), ('the', 'DT'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (',', 'VBP'), ('The', 'DT'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (',', 'NNP'), ('III', 'NN e', 'VBP'), ('the', 'DT'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (', ', ', ','), ('running', 'VBG'), ('the', 'DT'), ('new', 'JJ'), ('Jelly', 'NNP'), ('Bean', 'NN P'), ('system', 'NN'), (', ', ', '), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('8.9', 'CD'), ('Wifi', 'NNP'), ('tablet', 'NNP'), (', ', ', '), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (', ', ', '), ('Galaxy', 'NN P'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), ('mini', 'NN'), ('.', '.')], [('\n', '\_SP'), ('Apple', 'NNP'), ('state d', 'VBD'), ('iff', 'PRD'), ('had', 'VBD'), ('affeacted', 'VBN'), ('guickly', 'RB'), ('a d', 'VBD'), ('it', 'PRP'), ('had', 'VBD'), ('"acted', 'VBN'), ('quickly', 'RB'), ('a nd', 'CC'), ('diligently', 'RB'), ('"', "''"), ('in', 'IN'), ('order', 'NN'), ('to', 'TO'), ('"', '``'), ('determine', 'VB'), ('that', 'IN'), ('these', 'DT'), ('newly', 'RB'), ('that', 'IN'), ('these', 'DT'), ('t B'), ('released', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe' ny', 'JJ'), ('of', 'IN'), ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('alread y', 'RB'), ('asserted', 'VBN'), ('by', 'IN'), ('Apple', 'NNP'), ('.', '.'), ('"', "''")], [('\n', '\_SP'), ('In', 'IN'), ('August', 'NNP'), (',', ','), ('Samsung', 'NN P'), ('lost', 'VBD'), ('a', 'DT'), ('US', 'NNP'), ('patent', 'NN'), ('case', 'NN'), ('to', 'IN'), ('Apple', 'NNP'), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('t o', 'TO'), ('pay', 'VB'), ('its', 'PRP\$'), ('rival', 'JJ'), ('\$', '\$'), ('1.05bn', 'C D'), ('(', '-LRB-'), ('£0.66bn', 'NNP'), (')', '-RRB-'), ('in', 'IN'), ('damages', 'N NS'), ('for', 'IN'), ('copying', 'VBG'), ('features', 'NNS'), ('of', 'IN'), ('the', 'D NS'), ('for', 'IN'), ('copying', 'VBG'), ('features', 'NNS'), ('of', 'IN'), ('the', 'D T'), ('iPad', 'NNP'), ('and', 'CC'), ('iPhone', 'NNP'), ('in', 'IN'), ('its', 'PRP\$'), ('Galaxy', 'NNP'), ('range', 'NN'), ('of', 'IN'), ('devices', 'NNS'), ('.', '.')], [('Samsung', 'NNP'), (',', ','), ('which', 'WDT'), ('is', 'VBZ'), ('the', 'DT'), ('world', 'NN'), ("'s", 'POS'), ('top', 'JJ'), ('mobile', 'JJ'), ('phone', 'NN'), ('maker', 'NN'), (',', ','), ('is', 'VBZ'), ('appealing', 'VBG'), ('the', 'DT'), ('ruling', 'N N'), ('.', '.')], [('\n', '\_SP'), ('A', 'DT'), ('similar', 'JJ'), ('case', 'NN'), ('in', 'IN'), ('the', 'DT'), ('similar', 'JJ'), ('Samsung', 'NNP'), ("'s", 'POS'), ('favour', 'NNP'), ('found', 'VBD'), ('in', 'IN'), ('Apple', 'NNP'), ('to', 'TO'), ('publish', 'VB'), ('an', 'DT'), ('apology', 'NN'), ('making', 'VBG'), ('clear', 'JJ'), ('that', 'IN'), ('the', 'DT'), ('South', 'JJ'), ('Korean', 'JJ'), ('firm', 'NN'), ('had', 'VBD'), ('not', 'RB'), ('copied', 'VBN'), ('its', 'PRP) J'), ('firm', 'NN'), ('had', 'VBD'), ('not', 'RB'), ('copied', 'VBN'), ('its', 'PRP \$'), ('iPad', 'NNP'), ('when', 'WRB'), ('designing', 'VBG'), ('its', 'PRP\$'), ('own', 'JJ'), ('devices', 'NNS'), ('.', '.')]] In [24]: print(ner\_tags\_per\_sentence\_spacy) [[('San Jose', 'GPE'), ('California', 'GPE'), ('November 23', 'DATE'), ('six', 'CARDIN AL'), ('Samsung', 'ORG'), ('Jelly Bean', 'WORK\_OF\_ART'), ('Apple', 'ORG')], [], [('si x', 'CARDINAL'), ('the Galaxy S III', 'GPE'), ('Jelly Bean', 'ORG'), ('the Galaxy Tab 2 10.1', 'ORG')], [('Apple', 'ORG')], [('August', 'DATE'), ('Samsung', 'ORG'), ('US', 'GPE'), ('Apple', 'ORG'), ('1.05bn', 'MONEY'), ('iPad', 'ORG'), ('iPhone', 'ORG')], [('Samsung', 'ORG')], [('UK', 'GPE'), ('Samsung', 'ORG'), ('Apple', 'ORG'), ('South Ko rean', 'NORP'), ('iPad', 'ORG')]] In [97]: print(constituency\_output\_per\_sentence\_spacy)  $\hbox{\tt [[('https://www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more and of the context of the contex$ e-products-under-scrutiny.html', 'compound', ), ('\n\n', 'dep', ), ('Documents', 'appos', ), ('filed', 'acl', Documents), ('to', 'prep', filed), ('the', 'det', court), ('San', 'nmod', Jose), ('Jose', 'nmod', court), ('federal', 'amod', court), ('court', 'pobj', to), ('in', 'prep', court), ('California', 'pobj', in), ('on', 'prep', filed), ('Novem ber', 'pobj', on), ('23', 'nummod', November), ('list', 'appos', ), ('six', 'nummod', products), ('Samsung', 'compound', products), ('products', 'appo s', list), ('running', 'acl', products), ('the', 'det', Bean), ('"', 'punct', Bean), ('Jelly', 'compound', Bean', 'dobj', running), ('"', 'punct', Bean), ('and', 'cc', Bean), ('"', 'punct', Sandwich), ('Ice', 'compound', Cream), ('Cream', 'compoun d', Sandwich), ('Sandwich', 'conj', Bean), ('"', 'punct', Sandwich), ('operating', 'co mpound', systems), ('systems', 'appos', list), (',', 'punct', systems), ('which', 'nsu bj', infringe), ('Apple', 'compound', claims), ('claims', 'nsubj', infringe), ('infringe) ge', 'relcl', systems), ('its', 'poss', patents), ('patents', 'dobj', infringe), ('.', 'punct', )], [('\n', 'dep', )], [('The', 'det', phones), ('six', 'nummod', phones), ('phones', 'nsubj', are), ('an d', 'cc', phones), ('tablets', 'conj', phones), ('affected', 'acl', phones), ('are', 'ROOT', are), ('the', 'det', III), ('Galaxy', 'compound', III), ('S', 'compound', II I), ('III', 'attr', are), (',', 'punct', are), ('running', 'advcl', are), ('the', 'det', system), ('new', 'amod', system), ('Jelly', 'compound', Bean), ('Bean', 'compound') d', system), ('system', 'dobj', running), (',', 'punct', system), ('the', 'det', Tab), ('Galaxy', 'compound', Tab), ('Tab', 'appos', system), ('8.9', 'nummod', tablet), ('Wifi', 'compound', tablet), ('tablet', 'appos', Tab), (',', 'punct', Tab), ('the', 'det', Tab), ('Galaxy', 'compound', Tab), ('Tab', 'appos', Tab), ('2', 'nummod', Tab), ('10.1', 'nummod', Tab), (',', 'punct', system), ('Galaxy', 'compound', Rugby), ('Rugb y', 'appos', system), ('Pro', 'appos', Rugby), ('and', 'cc', Pro), ('Galaxy', 'conj', Pro), ('S', 'conj', Pro), ('III', 'nummod', mini), ('mini', 'attr', are), ('.', 'punc t', are)], [('\n', 'dep', stated), ('Apple', 'nsubj', stated), ('stated', 'ROOT', stated), ('it', 'nsubj', "acted), ('had', 'aux', "acted), ('"acted', 'ccomp', stated) d), ('quickly', 'advmod', "acted), ('and', 'cc', quickly), ('diligently', 'conj', qu ickly), ('"', 'punct', "acted), ('in', 'prep', "acted), ('order', 'pobj', in), ('t o', 'aux', determine), ('"', 'punct', determine), ('determine', 'acl', order), ('tha t', 'mark', infringe), ('these', 'det', products), ('newly', 'advmod', released), ('re leased', 'amod', products), ('products', 'nsubj', infringe), ('do', 'aux', infringe), ('infringe', 'ccomp', determine), ('many', 'dobj', infringe), ('of', 'prep', many), ('the', 'det', claims), ('same', 'amod', claims), ('claims', 'pobj', of), ('already', 'advmod', asserted), ('asserted', 'acl', claims), ('by', 'agent', asserted), ('Apple', 'pobj', by), ('.', 'punct', stated)], [('\n', 'dep', lost), ('Int', 'prep', lost), ('Sameung', 'psub', 'prep', lost), ('Sameung', 'psub', 'psub' ('In', 'prep', lost), ('August', 'pobj', In), (',', 'punct', lost), ('Samsung', 'nsub', lost), ('August', 'pobj', In'), (',', 'punct', lost), ('Samsung', 'nsub', lost), ('August', 'pobj', In'), (',', 'punct', lost), ('Samsung', 'nsub', lost), ('August', 'pobj', In'), (',', 'punct', lost), ('Samsung', 'nsub', lost), ('August', 'pobj', In'), (',', 'punct', lost), ('Samsung', 'nsub', 'nsub',j', lost), ('lost', 'ROOT', lost), ('a', 'det', case), ('US', 'compound', case), ('pat ent', 'compound', case), ('case', 'dobj', lost), ('to', 'prep', lost), ('Apple', 'pob j', to), ('and', 'cc', lost), ('was', 'auxpass', ordered), ('ordered', 'conj', lost), ('to', 'aux', pay), ('pay', 'xcomp', ordered), ('its', 'poss', 1.05bn), ('rival', 'amo d', 1.05bn), ('\$', 'nmod', 1.05bn), ('1.05bn', 'dobj', pay), ('(', 'punct', 1.05bn), ('1.05bn', 'dobj', pay), ('(', 'punct', 1.05bn), ('lost), 'lost), 'l ('£0.66bn', 'appos', 1.05bn), (')', 'punct', 1.05bn), ('in', 'prep', pay), ('damage s', 'pobj', in), ('for', 'prep', damages), ('copying', 'pcomp', for), ('features', 'do bj', copying), ('of', 'prep', features), ('the', 'det', iPad), ('iPad', 'pobj', of), ('and', 'cc', iPad), ('iPhone', 'conj', iPad), ('in', 'prep', copying), ('its', 'pos s', range), ('Galaxy', 'compound', range), ('range', 'pobj', in), ('of', 'prep', rang e), ('devices', 'pobj', of), ('.', 'punct', lost)], [('Samsung', 'nsubj', appealing), (',', 'punct', Samsung), ('which', 'nsubj', is), ('is', 'relcl', Samsung), ('the', 'de t', world), ('world', 'poss', maker), ("'s", 'case', world), ('top', 'amod', maker), ('mobile', 'amod', phone), ('phone', 'compound', maker), ('maker', 'attr', is), (',' 'punct', Samsung), ('is', 'aux', appealing), ('appealing', 'ROOT', appealing), ('the', 'det', ruling), ('ruling', 'dobj', appealing), ('.', 'punct', appealing)], [('\n', 'de p', case), ('A', 'det', case), ('similar', 'amod', case), ('case', 'ROOT', case), ('i n', 'prep', case), ('the', 'det', UK), ('UK', 'pobj', in), ('found', 'acl', case), ('in', 'prep', found), ('Samsung', 'poss', favour), ("'s", 'case', Samsung), ('favour', 'pobj', in), ('and', 'cc', found), ('ordered', 'conj', found), ('Apple', 'dobj', order ed), ('to', 'aux', publish), ('publish', 'xcomp', ordered), ('an', 'det', apology), ('apology', 'dobj', publish), ('making', 'acl', apology), ('clear', 'acomp', making), ('that', 'mark', copied), ('the', 'det', firm), ('South', 'amod', Korean), ('Korean', 'amod', firm), ('firm', 'nsubj', copied), ('had', 'aux', copied), ('not', 'neg', copie d), ('copied', 'ccomp', making), ('its', 'poss', iPad), ('iPad', 'dobj', copied), ('wh en', 'advmod', designing), ('designing', 'advcl', copied), ('its', 'poss', devices), ('own', 'amod', devices), ('devices', 'dobj', designing), ('.', 'punct', case)]] small tip: You can use **sents = list(doc.sents)** to be able to use the index to access a sentence like **sents[2]** for the third sentence. [total points: 7] Exercise 3: Comparison NLTK and spaCy We will now compare the output of NLTK and spaCy, i.e., in what do they differ? [points: 3] Exercise 3a: Part of speech tagging Compare the output from NLTK and spaCy regarding part of speech tagging. To compare, you probably would like to compare sentence per sentence. Describe if the sentence splitting is different for NLTK than for spaCy. If not, where do they differ? After checking the sentence splitting, select a sentence for which you expect interesting results and perhaps differences. Motivate your choice. Compare the output in token.tag from spaCy to the part of speech tagging from NLTK for each token in your selected sentence. Print for each token the output from NLTK and spaCy next to each other (align possible tokenization differences). Are there any differences? This is not a trick question; it is possible that there are no differences. #Comparing sentence splitting for NLTK and for spaCy. In [48]: print(pos tags per sentence,'\n') print(pos tags per sentence spacy) [[('https', ('//www.telegraph.co.uk/technology/apple/9702716/Apple Samsung-lawsuit-six-more-products-under-scrutiny.html', 'JJ'), ('Documents', 'NNS'), ('filed', 'VBN'), ('to', 'TO'), ('the', 'DT'), ('San', 'NNP'), ('Jose', 'NNP'), ('fede ral', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Nov ember', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('pro ducts', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('``', '``'), ('Jelly', 'RB'), ('Be an', 'NNP'), ("''", "''"), ('and', 'CC'), ('``', '``'), ('Ice', 'NNP'), ('Cream', 'NN P'), ('Sandwich', 'NNP'), (""'", ""'"), ('operating', 'VBG'), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), ('claims', 'VBZ'), ('infringe', 'VB'), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')], [('The', 'DT'), ('six', 'CD'), ('pho nes', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), ('affected', 'VBN'), ('are', 'VBP') nes', 'NNS'), ('and', 'CC'), ('tablets', 'NNS'), ('affected', 'VBN'), ('are', 'VBP'), ('the', 'DT'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (',', ','), ('runnin g', 'VBG'), ('the', 'DT'), ('new', 'JJ'), ('Jelly', 'NNP'), ('Bean', 'NNP'), ('syste m', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('8.9', 'C D'), ('Wifi', 'NNP'), ('tablet', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), ('Galaxy', 'NNP'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NN P'), ('mini', 'NN'), ('.', '.')], [('Apple', 'NNP'), ('stated', 'VBD'), ('it', 'PRP'), ('had', 'VBD'), ('"acted', 'VBN'), ('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ("''", "''"), ('in', 'IN'), ('order', 'NN'), ('to', 'TO'), ('``', '``'), ('dete rmine', 'VB'), ('that', 'IN'), ('these', 'DT'), ('newly', 'RB'), ('released', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'JJ'), ('of', 'IN'), ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('already', 'RB'), ('asserted', 'VB N'), ('by', 'IN'), ('Apple', 'NNP'), ('.', '.'), ("''", "''")], [('In', 'IN'), ('Augus t', 'NNP'), (',', ','), ('Samsung', 'NNP'), ('lost', 'VBD'), ('a', 'DT'), ('US', 'NN P'), ('patent', 'NN'), ('case', 'NN'), ('to', 'TO'), ('Apple', 'NNP'), ('and', 'CC'), P'), ('patent', 'NN'), ('case', 'NN'), ('to', 'TO'), ('Apple', 'NNP'), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('to', 'TO'), ('pay', 'VB'), ('its', 'PRP\$'), ('rival', 'JJ'), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '('), ('£0.66bn', 'NN'), (')', ')'), ('in', 'IN'), ('damages', 'NNS'), ('for', 'IN'), ('copying', 'VBG'), ('features', 'NNS'), ('of', 'IN'), ('the', 'DT'), ('ipad', 'NN'), ('and', 'CC'), ('iphone', 'NN'), ('in', 'IN'), ('its', 'PRP\$'), ('Galaxy', 'NNP'), ('range', 'NN'), ('of', 'IN'), ('devices', 'NNS'), ('.', '.')], [('Samsung', 'NNP'), (', ', ','), ('which', 'WDT'), ('is', 'VBZ'), ('the', 'DT'), ('world', 'NN'), ("'s", 'POS'), ('top', 'JJ'), ('mobile', 'NN'), ('phone', 'NN'), ('maker', 'NN'), (', ', ', '), ('is', 'VBZ'), ('appealing', 'VBG'), ('the', 'DT'), ('the', 'DT'), ('similar', 'JJ'), ('case', 'NN'), ('in', 'IN'), ('the', 'DT'), ('found', 'VBD'), ('in', 'IN'), ('samsung', 'NNP'), ("s", 'POS'), ('favour', 'NNP'), ('found', 'VBD'), ('in', 'IN'), ('making', 'VBG'), ('clear', 'JJ'), ('that', 'IN'), ('the', 'DT'), ('south', 'JJ'), ('korean', 'JJ'), ('firm', 'NN'), ('had', 'VBD'), ('not', 'RB'), ('copied', 'VBN'), ('its', 'PRP\$'), ('ipad', 'NN'), ('when', 'WRB'), ('designing', 'VBG'), ('its', 'PRP\$'), ('devices', 'NNS'), ('.', '.')]] [[('https://www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-mor e-products-under-scrutiny.html', 'NNP'), ('\n\n', '\_SP'), ('Documents', 'NNS'), ('file d', 'VBD'), ('to', 'IN'), ('the', 'DT'), ('San', 'NNP'), ('Jose', 'NNP'), ('federal', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Novembe r', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('product s', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('"', '``'), ('Jelly', 'NNP'), ('Bean', 'NNP'), ('"', "''"), ('and', 'CC'), ('"', '``'), ('Ice', 'NNP'), ('Cream', 'NNP'), ('S andwich', 'NNP'), ('"', "''"), ('operating', 'NN'), ('systems', 'NNS'), (',', ','), e', 'VBP'), ('the', 'DT'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), (',', ','), ('running', 'VBG'), ('the', 'DT'), ('new', 'JJ'), ('Jelly', 'NNP'), ('Bean', 'NN P'), ('system', 'NN'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('8.9', 'CD'), ('Wifi', 'NNP'), ('tablet', 'NNP'), (',', ','), ('the', 'DT'), ('Galaxy', 'NNP'), ('Tab', 'NNP'), ('2', 'CD'), ('10.1', 'CD'), (',', ','), ('Galaxy', 'NN P'), ('Rugby', 'NNP'), ('Pro', 'NNP'), ('and', 'CC'), ('Galaxy', 'NNP'), ('S', 'NNP'), ('III', 'NNP'), ('mini', 'NN'), ('.', '.')], [('\n', '\_SP'), ('Apple', 'NNP'), ('state d', 'VBD'), ('it', 'PRP'), ('had', 'VBD'), ('â&acted', 'VBN'), ('quickly', 'RB'), ('and', 'CC'), ('diligently', 'RB'), ('"', ""'"), ('in', 'IN'), ('order', 'NN'), ('to', 'TO'), ('"', '``), ('determine', 'VB'), ('that', 'IN'), ('these', 'DT'), ('newly', 'RB'), ('released', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('ma B'), ('released', 'VBN'), ('products', 'NNS'), ('do', 'VBP'), ('infringe', 'VB'), ('many', 'JJ'), ('of', 'IN'), ('the', 'DT'), ('same', 'JJ'), ('claims', 'NNS'), ('already', 'RB'), ('asserted', 'VBN'), ('by', 'IN'), ('Apple', 'NNP'), ('.', '.'), ('"', "''")], [('\n', '\_SP'), ('In', 'IN'), ('August', 'NNP'), (',', ','), ('Samsung', 'NNP'), ('lost', 'VBD'), ('a', 'DT'), ('US', 'NNP'), ('patent', 'NN'), ('case', 'NN'), ('the stand', 'NN'), ('the sta ('to', 'IN'), ('Apple', 'NNP'), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('t ('to', 'IN'), ('Apple', 'NNP'), ('and', 'CC'), ('was', 'VBD'), ('ordered', 'VBN'), ('to', 'TO'), ('pay', 'VB'), ('its', 'PRP\$'), ('rival', 'JJ'), ('\$', '\$'), ('1.05bn', 'CD'), ('(', '-LRB-'), ('\hat{A}\tilde{0.66bn'}, 'NNP'), (')', '-RRB-'), ('in', 'IN'), ('damages', 'NNS'), ('for', 'IN'), ('copying', 'VBG'), ('features', 'NNS'), ('of', 'IN'), ('the', 'DT'), ('iPad', 'NNP'), ('and', 'CC'), ('iPhone', 'NNP'), ('in', 'IN'), ('its', 'PRP\$'), ('Galaxy', 'NNP'), ('range', 'NN'), ('of', 'IN'), ('devices', 'NNS'), ('.', '.')], [('Samsung', 'NNP'), (',', ','), ('which', 'WDT'), ('is', 'VBZ'), ('the', 'DT'), ('world', 'NN'), ("'s", 'POS'), ('top', 'JJ'), ('mobile', 'JJ'), ('phone', 'NN'), ('maker', 'NN'), (',', ','), ('is', 'VBZ'), ('appealing', 'VBG'), ('the', 'DT'), ('ruling', 'NN'), ('.', '.')], [('\n', '\_SP'), ('A', 'DT'), ('similar', 'JJ'), ('case', 'NN'), ('in', 'IN'), ('the', 'DT'), ('UK', 'NNP'), ('found', 'VBD'), ('in', 'IN'), ('Samsung'. n', 'IN'), ('the', 'DT'), ('UK', 'NNP'), ('found', 'VBD'), ('in', 'IN'), ('Samsung', 'NNP'), ("'s", 'POS'), ('favour', 'NN'), ('and', 'CC'), ('ordered', 'VBD'), ('Apple', 'NNP'), ('to', 'TO'), ('publish', 'VB'), ('an', 'DT'), ('apology', 'NN'), ('making', 'VBG'), ('clear', 'JJ'), ('that', 'IN'), ('the', 'DT'), ('South', 'JJ'), ('Korean', 'J J'), ('firm', 'NN'), ('had', 'VBD'), ('not', 'RB'), ('copied', 'VBN'), ('its', 'PRP \$'), ('iPad', 'NNP'), ('when', 'WRB'), ('designing', 'VBG'), ('its', 'PRP\$'), ('own', 'JJ'), ('devices', 'NNS'), ('.', '.')]] In [41]: | #We can see that there are indeed some differences for nltk and spacy. #For example, in one of the first sentences the word 'to' is labeled as 'TO' with NLT #we can also observe that some of the splitting is different. In [54]: | print(pos\_tags\_per\_sentence[0], '\n') print(pos tags per sentence spacy[0]) #We choose the sentence below because we are interested to see how the website url is #We think that NLTK and spacy might split the url in a different way. #In addition, we think that some of the words will be labeled differently. [('https', 'NN'), (':', ':'), ('//www.telegraph.co.uk/technology/apple/9702716/Apple-S amsung-lawsuit-six-more-products-under-scrutiny.html', 'JJ'), ('Documents', 'NNS'), ('filed', 'VBN'), ('to', 'TO'), ('the', 'DT'), ('San', 'NNP'), ('Jose', 'NNP'), ('fede ral', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Nov ember', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('products', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('`', '``'), ('Jelly', 'RB'), ('Be an', 'NNP'), ("''", "''"), ('and', 'CC'), ('``', '``'), ('Ice', 'NNP'), ('Cream', 'NN P'), ('Sandwich', 'NNP'), ("''", "''"), ('operating', 'VBG'), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), ('claims', 'VBZ'), ('infringe', 'VB'), ('its', 'PRP\$'), ('patents', 'NNS'), ('.', '.')] [('https://www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more -products-under-scrutiny.html', 'NNP'), ('\n\n', '\_SP'), ('Documents', 'NNS'), ('file d', 'VBD'), ('to', 'IN'), ('the', 'DT'), ('San', 'NNP'), ('Jose', 'NNP'), ('federal', 'JJ'), ('court', 'NN'), ('in', 'IN'), ('California', 'NNP'), ('on', 'IN'), ('Novembe r', 'NNP'), ('23', 'CD'), ('list', 'NN'), ('six', 'CD'), ('Samsung', 'NNP'), ('product s', 'NNS'), ('running', 'VBG'), ('the', 'DT'), ('"', '``'), ('Jelly', 'NNP'), ('Bean', 'NNP'), ('"', "''"), ('and', 'CC'), ('"', '``'), ('Ice', 'NNP'), ('Cream', 'NNP'), ('s andwich', 'NNP'), ('"', "''"), ('operating', 'NN'), ('systems', 'NNS'), (',', ','), ('which', 'WDT'), ('Apple', 'NNP'), ('claims', 'NNS'), ('infringe', 'VBP'), ('its', 'PDE'), ('patents', 'NNS'), (',', ','), RP\$'), ('patents', 'NNS'), ('.', '. **for** item a, item b in zip(pos tags per sentence[0], pos tags per sentence spacy[0]): print(item a, item b) ('https', 'NN') ('https://www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-l awsuit-six-more-products-under-scrutiny.html', 'NNP') (':', ':') ('\n\n', 'SP') ('//www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more-produc ts-under-scrutiny.html', 'JJ') ('Documents', 'NNS') ('Documents', 'NNS') ('filed', 'VBD') ('filed', 'VBN') ('to', 'IN') ('to', 'TO') ('the', 'DT') ('the', 'DT') ('San', 'NNP') ('San', 'NNP') ('Jose', 'NNP') ('Jose', 'NNP') ('federal', 'JJ') ('federal', 'JJ') ('court', 'NN') ('court', 'NN') ('in', 'IN') ('in', 'IN') ('California', 'NNP') ('California', 'NNP') ('on', 'IN') ('on', 'IN') ('November', 'NNP') ('November', 'NNP') ('23', 'CD') ('23', 'CD') ('list', 'NN') ('list', 'NN') ('six', 'CD') ('six', 'CD') ('Samsung', 'NNP') ('Samsung', 'NNP') ('products', 'NNS') ('products', 'NNS') ('running', 'VBG') ('running', 'VBG') ('the', 'DT') ('the', 'DT') ('"', '``') ('``', '``') ('Jelly', 'NNP') ('Jelly', 'RB') ('Bean', 'NNP') ('Bean', 'NNP') ('"', "''") ("'', "''") ('and', 'CC') ('"', '``') ('``', '``') ('Ice', 'NNP') ('Ice', 'NNP') ('Cream', 'NNP') ('Cream', 'NNP') ('Sandwich', 'NNP') ('Sandwich', 'NNP') ('"', "''") ("''", "''") ('operating', 'NN') ('operating', 'VBG') ('systems', 'NNS') 'NNS') (', (',', ',') ('which', 'WDT') ('which', 'WDT') ('Apple', 'NNP') ('Apple', 'NNP') ('claims', 'NNS') ('claims', 'VBZ') ('infringe', 'VBP') ('infringe', 'VB') ('its', 'PRP\$')
('its', 'PRP\$') ('patents', 'NNS') ('patents', 'NNS') ('.', '.') In [ ]: # In the sentences above we can observe some differences between the behavior in NLTK #First, we can see that not every words is split the same way. The URL at the beginning #while spacy retains the whole URL. Secondly, #we can see that the words 'jelly', 'operating' and the word 'claims' are labeled dif. [points: 2] Exercise 3b: Named Entity Recognition (NER) Describe differences between the output from NLTK and spaCy for Named Entity Recognition. Which one do you think performs better? In [64]: for i in ner tags per sentence: print(i) for i in ner\_tags\_per\_sentence\_spacy: print(i) (S https/NN //www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-more-produc ts-under-scrutiny.html/JJ Documents/NNS filed/VBN to/TO the/DT (ORGANIZATION San/NNP Jose/NNP) federal/JJ court/NN in/IN (GPE California/NNP) on/IN November/NNP 23/CD list/NN six/CD (ORGANIZATION Samsung/NNP) products/NNS running/VBG the/DT Jelly/RB (GPE Bean/NNP) 11/11 and/CC `/` Ice/NNP Cream/NNP Sandwich/NNP 11/11 operating/VBG systems/NNS ,/, which/WDT (PERSON Apple/NNP) claims/VBZ infringe/VB its/PRP\$ patents/NNS ./.) (S The/DT six/CD phones/NNS and/CC tablets/NNS affected/VBN are/VBP the/DT (ORGANIZATION Galaxy/NNP) S/NNP III/NNP running/VBG the/DT new/JJ (PERSON Jelly/NNP Bean/NNP) system/NN the/DT (ORGANIZATION Galaxy/NNP) Tab/NNP 8.9/CD Wifi/NNP tablet/NN ,/, the/DT (ORGANIZATION Galaxy/NNP) 2/CD 10.1/CD (PERSON Galaxy/NNP Rugby/NNP Pro/NNP) and/CC (PERSON Galaxy/NNP S/NNP) III/NNP mini/NN (S (PERSON Apple/NNP) stated/VBD it/PRP had/VBD "acted/VBN quickly/RB and/CC diligently/RB 11/11 in/IN order/NN to/TO determine/VB that/IN these/DT newly/RB released/VBN products/NNS do/VBP infringe/VB many/JJ of/IN the/DT same/JJ claims/NNS already/RB asserted/VBN by/IN (PERSON Apple/NNP) (S In/IN (GPE August/NNP) (PERSON Samsung/NNP) lost/VBD a/DT (GSP US/NNP) patent/NN case/NN to/TO (GPE Apple/NNP) and/CC was/VBD ordered/VBN to/TO pay/VB its/PRP\$ rival/JJ

	)/) in/IN damages/NNS for/IN copying/VBG features/NNS of/IN the/DT (ORGANIZATION iPad/NN) and/CC (ORGANIZATION iPhone/NN) in/IN its/PRP\$ (GPE Galaxy/NNP) range/NN of/IN devices/NNS ./.) (S (GPE Samsung/NNP) ,/, which/WDT is/VBZ the/DT world/NN 's/POS top/JJ pebile/NN
	<pre>top/JJ mobile/NN phone/NN maker/NN ,/, is/VBZ appealing/VBG the/DT ruling/NN ./.) (S     A/DT     similar/JJ     case/NN     in/IN     the/DT     (ORGANIZATION UK/NNP)     found/VBD     in/IN     (GPE Samsung/NNP) 's/POS favour/NN and/CC ordered/VBD</pre>
	and/CC
[65]:	<pre>its/PRP\$ own/JJ devices/NNS ./.) [('San Jose', 'GPE'), ('California', 'GPE'), ('November 23', 'DATE'), ('six', 'CARI L'), ('Samsung', 'ORG'), ('Jelly Bean', 'WORK_OF_ART'), ('Apple', 'ORG')] [] [('six', 'CARDINAL'), ('the Galaxy S III', 'GPE'), ('Jelly Bean', 'ORG'), ('the Gal Tab 2 10.1', 'ORG')] [('Apple', 'ORG')] [('Apple', 'ORG')] [('August', 'DATE'), ('Samsung', 'ORG'), ('US', 'GPE'), ('Apple', 'ORG'), ('1.05bn') 'MONEY'), ('iPad', 'ORG'), ('iPhone', 'ORG')] [('Samsung', 'ORG')] [('UK', 'GPE'), ('Samsung', 'ORG'), ('Apple', 'ORG'), ('South Korean', 'NORP'), ('id', 'ORG')]  # SpaCy finds the relations betweeen words instead of seeing all the words as sepal # a word like 'jelly bean' is split into 'jelly' and 'bean' in NLTK and spacy class')</pre>
[68]:	# Another example is the word 'the Galaxy S III' which is seen as an individual word # Eventhough spacy incorrectly labeled 'the Galaxy S III' as 'GPE' meaning Geopolis # important distinction that some words that should be seen as one entity are correctly and the space of the seen as one entity are correctly and the space of the seen as one entity are correctly and the seen as one entity are correctly are seen as one entity are correctly are correctly are seen as one entity are correctly are correctly are seen as one entity are correctly are co
	<pre>(S     (NP https/NN) :/:     (NP</pre>
	<pre>(NP list/NN) six/CD Samsung/NNP products/NNS (VP (V running/VBG) (NP the/DT)) ``/` Jelly/RB Bean/NNP ''/' and/CC ``/` Ice/NNP Cream/NNP Sandwich/NNP ''/'' (VP (V operating/VBG)) systems/NNS ,/, which/WDT</pre>
	<pre>which/WDT Apple/NNP (VP (V claims/VBZ)) (VP (V infringe/VB)) its/PRP\$ patents/NNS ./.)  [('https://www.telegraph.co.uk/technology/apple/9702716/Apple-Samsung-lawsuit-six-r-products-under-scrutiny.html', 'compound', ), ('\n\n', 'dep', ), ('Documents', 'appos', ), ('filed', 'acl', Documents), ('to', 'prep', filed), ('the', 'det', court), ('Sam', 'appos', 'Jose', 'nmod', court), ('federal', 'amod', court), ('court', 'pobjto), ('in', 'prep', court), ('California', 'pobj', in), ('on', 'prep', filed), ('Nober', 'pobj', on), ('23', 'nummod', November), ('list', 'appos',</pre>
ı [ ]:	), ('six', 'nummod', products), ('Samsung', 'compound', products), ('products', 'ars', list), ('running', 'acl', products), ('the', 'det', Bean), ('"', 'punct', Bean) ('Jelly', 'compound', Bean), ('Bean', 'dobj', running), ('"', 'punct', Bean), ('and 'cc', Bean), ('"', 'punct', Sandwich), ('Ice', 'compound', Cream), ('Cream', 'compound', Sandwich), ('Sandwich', 'conj', Bean), ('"', 'punct', Sandwich), ('operating', mpound', systems), ('systems', 'appos', list), (',', 'punct', systems), ('which', 'bj', infringe), ('Apple', 'compound', claims), ('claims', 'nsubj', infringe), ('infige', 'relcl', systems), ('its', 'poss', patents), ('patents', 'dobj', infringe), ('punct',  )]  # Dependency parsing will show the relationships between words and their constituted while constituency parsing will showcase the the whole sentence structure and the shift is a library that processes strings. It takes strings as input and it returned in the shift is a library that processes strings. It takes strings as input and it returned in the shift is parsed with the shift is pars
	#Looking at the above output, we can observe differences between NLTK and spaCy. # For example, the word 'cream' is labeled NNP in NLTK while in spaCy it is labeled # some words seem to be labeled the same, but overall there are many differences to the same words notebook  End of this notebook