

225229110

HARI PRASATH S

LAB9:BUILDING BIGRAM TAGGER

```
In [5]: import nltk
nltk.download('averaged_perceptron_tagger')
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] C:\Users\1mscdsa10\AppData\Roaming\nltk_data...
[nltk_data] Unzipping taggers\averaged_perceptron_tagger.zip.
```

Out[5]: True

```
In [6]: from nltk.tokenize import *
```

```
In [7]: import nltk
text=word_tokenize('And now for something completely different')
nltk.pos_tag(text)
```

Out[7]:

```
[('And', 'CC'),
 ('now', 'RB'),
 ('for', 'IN'),
 ('something', 'NN'),
 ('completely', 'RB'),
 ('different', 'JJ')]
```

Ex-2

```
In [8]: import nltk
nltk.download('brown')
```

```
[nltk_data] Downloading package brown to
[nltk_data] C:\Users\1mscdsa10\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\brown.zip.
```

Out[8]: True

```
In [11]: from nltk.corpus import brown
tagsen=brown.tagged_sents()
```

In [12]: tagsen

Out[12]: [[('The', 'AT'), ('Fulton', 'NP-TL'), ('County', 'NN-TL'), ('Grand', 'JJ-TL'), ('Jury', 'NN-TL'), ('said', 'VBD'), ('Friday', 'NR'), ('an', 'AT'), ('investigation', 'NN'), ('of', 'IN'), ('Atlanta's', 'NP\$'), ('recent', 'JJ'), ('primary', 'NN'), ('election', 'NN'), ('produced', 'VBD'), ('`', '`'), ('no', 'AT'), ('evidence', 'NN'), ('"', '"'), ('that', 'CS'), ('any', 'DTI'), ('irregularities', 'NNS'), ('took', 'VBD'), ('place', 'NN'), ('.', '.')],
[('The', 'AT'), ('jury', 'NN'), ('further', 'RBR'), ('said', 'VBD'), ('in', 'IN'), ('term-end', 'NN'), ('presentments', 'NNS'), ('that', 'CS'), ('the', 'AT'), ('City', 'NN-TL'), ('Executive', 'JJ-TL'), ('Committee', 'NN-TL'), ('.', '.'), ('which', 'WDT'), ('had', 'HVD'), ('over-all', 'JJ'), ('charge', 'NN'), ('of', 'IN'), ('the', 'AT'), ('election', 'NN'), ('.', '.'), ('`', '`'), ('deserves', 'VBZ'), ('the', 'AT'), ('praise', 'NN'), ('and', 'CC'), ('thanks', 'NNS'), ('of', 'IN'), ('the', 'AT'), ('City', 'NN-TL'), ('of', 'IN-TL'), ('Atlanta', 'NP-TL'), ('"', '"'), ('for', 'IN'), ('the', 'AT'), ('manner', 'NN'), ('in', 'IN'), ('which', 'WDT'), ('the', 'AT'), ('election', 'NN'), ('was', 'BEDZ'), ('conducted', 'VBN'), ('.', '.')], ...]

In [14]: `import nltk`
`nltk.download('universal_tagset')`

[nltk_data] Downloading package universal_tagset to
[nltk_data] C:\Users\1mscdsa10\AppData\Roaming\nltk_data...
[nltk_data] Unzipping taggers\universal_tagset.zip.

Out[14]: True

In [15]: `from nltk.corpus import brown`
`brown_news_tagged = brown.tagged_sents(categories='news', tagset='universal')`
`brown_news_words = brown.tagged_words(categories='news', tagset='universal')`

`brown_train = brown_news_tagged[57340:]`
`brown_test = brown_news_tagged[:50000]`

`from nltk.tag import untag`
`test_sent = untag(brown_test[0])`
`print("Tagged: ", brown_test[0])`
`print("Untagged: ", test_sent)`

Tagged: [('The', 'DET'), ('Fulton', 'NOUN'), ('County', 'NOUN'), ('Grand', 'ADJ'), ('Jury', 'NOUN'), ('said', 'VERB'), ('Friday', 'NOUN'), ('an', 'DET'), ('investigation', 'NOUN'), ('of', 'ADP'), ('Atlanta's', 'NOUN'), ('recent', 'ADJ'), ('primary', 'NOUN'), ('election', 'NOUN'), ('produced', 'VERB'), ('`', '.'), ('no', 'DET'), ('evidence', 'NOUN'), ('"', '.'), ('that', 'ADP'), ('any', 'DET'), ('irregularities', 'NOUN'), ('took', 'VERB'), ('place', 'NOUN'), ('.', '.')]
Untagged: ['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', 'Friday', 'an', 'investigation', 'of', 'Atlanta's', 'recent', 'primary', 'election', 'produced', '`', 'no', 'evidence', '"', 'that', 'any', 'irregularities', 'took', 'place', '.']

```
In [18]: br_train = tagsen[0:50000]
br_test = tagsen[50000:]
br_test[0]
```

```
Out[18]: [('I', 'PPSS'),
('was', 'BEDZ'),
('loaded', 'VBN'),
('with', 'IN'),
('suds', 'NNS'),
('when', 'WRB'),
('I', 'PPSS'),
('ran', 'VBD'),
('away', 'RB'),
(',', ', '),
('and', 'CC'),
('I', 'PPSS'),
('haven't', 'HV*'),
('had', 'HVN'),
('a', 'AT'),
('chance', 'NN'),
('to', 'TO'),
('wash', 'VB'),
('it', 'PPO'),
('off', 'RP'),
('.', '. ')]
```

Step2:Build a bigram tagger

```
In [20]: t0 = nltk.DefaultTagger('NN')
t1 = nltk.UnigramTagger(br_train, backoff=t0)
t2 = nltk.BigramTagger(br_train, backoff=t1)
```

Step3:Evaluate

```
In [21]: t2.evaluate(br_test)
```

```
Out[21]: 0.9111100662708622
```

Step4:Explore

```
In [22]: # 1.
total_train = [len(l) for l in br_train]
sum(total_train)
```

```
Out[22]: 1039920
```

```
In [23]: total_test = [len(l) for l in br_test]
sum(total_test)
```

```
Out[23]: 121272
```

```
In [24]: # 2.  
t1.evaluate(br_test)
```

```
Out[24]: 0.8897849462365591
```

```
In [25]: t2.evaluate(br_test)
```

```
Out[25]: 0.9111006662708622
```

```
In [26]: # 3.  
br_train[0]
```

```
Out[26]: [('The', 'AT'),  
          ('Fulton', 'NP-TL'),  
          ('County', 'NN-TL'),  
          ('Grand', 'JJ-TL'),  
          ('Jury', 'NN-TL'),  
          ('said', 'VBD'),  
          ('Friday', 'NR'),  
          ('an', 'AT'),  
          ('investigation', 'NN'),  
          ('of', 'IN'),  
          ('Atlanta's', 'NP$'),  
          ('recent', 'JJ'),  
          ('primary', 'NN'),  
          ('election', 'NN'),  
          ('produced', 'VBD'),  
          ('', ''),  
          ('no', 'AT'),  
          ('evidence', 'NN'),  
          ('', ''),  
          ('that', 'CS'),  
          ('any', 'DTI'),  
          ('irregularities', 'NNS'),  
          ('took', 'VBD'),  
          ('place', 'NN'),  
          ('.', '.')] 
```

```
In [27]: br_train[1277]
```

```
Out[27]: [('', ''),  
          ('I', 'PPSS'),  
          ('told', 'VBD'),  
          ('him', 'PPO'),  
          ('who', 'WPS'),  
          ('I', 'PPSS'),  
          ('was', 'BEDZ'),  
          ('and', 'CC'),  
          ('he', 'PPS'),  
          ('was', 'BEDZ'),  
          ('quite', 'QL'),  
          ('cold', 'JJ'),  
          ('.', '.')] 
```

```
In [28]: br_train[1277] [11]
```

```
Out[28]: ('cold', 'JJ')
```

```
In [30]: # 4.  
br_train_flat = [(word, tag) for sent in br_train for (word, tag) in sent]
```

```
In [31]: br_train_flat[:40]
```

```
Out[31]: [('The', 'AT'),  
          ('Fulton', 'NP-TL'),  
          ('County', 'NN-TL'),  
          ('Grand', 'JJ-TL'),  
          ('Jury', 'NN-TL'),  
          ('said', 'VBD'),  
          ('Friday', 'NR'),  
          ('an', 'AT'),  
          ('investigation', 'NN'),  
          ('of', 'IN'),  
          ("Atlanta's", 'NP$'),  
          ('recent', 'JJ'),  
          ('primary', 'NN'),  
          ('election', 'NN'),  
          ('produced', 'VBD'),  
          ('', ''),  
          ('no', 'AT'),  
          ('evidence', 'NN'),  
          ('', ''),  
          ('that', 'CS'),  
          ('any', 'DTI'),  
          ('irregularities', 'NNS'),  
          ('took', 'VBD'),  
          ('place', 'NN'),  
          ('.', '.'),  
          ('The', 'AT'),  
          ('jury', 'NN'),  
          ('further', 'RBR'),  
          ('said', 'VBD'),  
          ('in', 'IN'),  
          ('term-end', 'NN'),  
          ('presentments', 'NNS'),  
          ('that', 'CS'),  
          ('the', 'AT'),  
          ('City', 'NN-TL'),  
          ('Executive', 'JJ-TL'),  
          ('Committee', 'NN-TL'),  
          ('', '''),  
          ('which', 'WDT'),  
          ('had', 'HVD')]
```

```
In [32]: br_train_flat[13]
```

```
Out[32]: ('election', 'NN')
```

```
In [33]: # 5. a)
fd = nltk.FreqDist(br_train_flat)
cfd = nltk.ConditionalFreqDist(br_train_flat)
```

```
In [34]: cfd['cold'].most_common()
```

```
Out[34]: [('JJ', 110), ('NN', 8), ('RB', 2)]
```

```
In [35]: # 5. b)
br_train_2grams = list(nltk.ngrams(br_train_flat, 2))
br_train_cold = [a[1] for (a,b) in br_train_2grams if b[0] == 'cold']
fdist = nltk.FreqDist(br_train_cold)
[tag for (tag, _) in fdist.most_common()]
```

```
Out[35]: ['AT',
          'IN',
          'CC',
          'QL',
          'BEDZ',
          'JJ',
          ',',
          'DT',
          'PP$',
          'RP',
          '\'',
          'NN',
          'VBN',
          'VBD',
          'CS',
          'BEZ',
          'DOZ',
          'RB',
          'PPSS',
          'BE',
          'VB',
          'VBZ',
          'NP$',
          'BEDZ*',
          '--',
          'DTI',
          'WRB',
          'BED']
```

```
In [36]: # 5. c)
br_pre = [(w2+" "+t2, t1) for ((w1,t1),(w2,t2)) in br_train_2grams]
br_pre_cfd = nltk.ConditionalFreqDist(br_pre)
br_pre
('Grand/JJ-TL', 'NN-TL'),
('Jury/NN-TL', 'JJ-TL'),
('said/VBD', 'NN-TL'),
('Friday/NR', 'VBD'),
('an/AT', 'NR'),
('investigation/NN', 'AT'),
('of/IN', 'NN'),
('Atlanta's/NP$', 'IN'),
('recent/JJ', 'NP$'),
('primary/NN', 'JJ'),
('election/NN', 'NN'),
('produced/VBD', 'NN'),
('``/``', 'VBD'),
('no/AT', '``'),
('evidence/NN', 'AT'),
(''''/'', 'NN'),
('that/CS', '''),
('any/DTI', 'CS'),
('irregularities/NNS', 'DTI'),
('took/VBD', 'NNS'),
...

```

```
In [37]: # 5. d)
br_pre_cfd['cold/NN'].most_common()
```

```
Out[37]: [('AT', 4), ('JJ', 2), ('', 1), ('DT', 1)]
```

```
In [38]: br_pre_cfd['cold/JJ'].most_common()
```

```
Out[38]: [('AT', 38),
          ('IN', 14),
          ('CC', 8),
          ('QL', 7),
          ('BEDZ', 7),
          ('JJ', 4),
          ('DT', 3),
          ('', 3),
          ('PP$', 3),
          ('``', 2),
          ('NN', 2),
          ('VBN', 2),
          ('VBD', 2),
          ('CS', 1),
          ('BEZ', 1),
          ('DOZ', 1),
          ('RB', 1),
          ('PPSS', 1),
          ('BE', 1),
          ('VB', 1),
          ('VBZ', 1),
          ('NP$', 1),
          ('BEDZ*', 1),
          ('--', 1),
          ('RP', 1),
          ('DTI', 1),
          ('WRB', 1),
          ('BED', 1)]
```

```
In [39]: # 6.
         bigram_tagger = nltk.BigramTagger(br_train)
```

```
In [40]: # 6. a)
         text1 = word_tokenize('I was very cold.')
         bigram_tagger.tag(text1)
```

```
Out[40]: [('I', 'PPSS'), ('was', 'BEDZ'), ('very', 'QL'), ('cold', 'JJ'), ('.', '.')]
```

```
In [41]: # 6. b)
         text2 = word_tokenize('I had a cold.')
         bigram_tagger.tag(text2)
```

```
Out[41]: [('I', 'PPSS'), ('had', 'HVD'), ('a', 'AT'), ('cold', 'JJ'), ('.', '.')]
```



```
In [42]: # 6. c)
text3 = word_tokenize('I had a severe cold.')
bigram_tagger.tag(text3)
```

```
Out[42]: [('I', 'PPSS'),
          ('had', 'HVD'),
          ('a', 'AT'),
          ('severe', 'JJ'),
          ('cold', 'JJ'),
          ('.', '.')]

```

```
In [43]: # 6. d)
text4 = word_tokenize('January was a cold month.')
bigram_tagger.tag(text4)
```

```
Out[43]: [('January', None),
          ('was', None),
          ('a', None),
          ('cold', None),
          ('month', None),
          ('.', None)]

```

```
In [44]: # 8. a)
text5 = word_tokenize('I failed to do so.')
bigram_tagger.tag(text5)
```

```
Out[44]: [('I', 'PPSS'),
          ('failed', 'VBD'),
          ('to', 'TO'),
          ('do', 'DO'),
          ('so', 'RB'),
          ('.', '.')]

```

```
In [45]: # 8. b)
text6 = word_tokenize('I was happy, but so was my enemy.')
bigram_tagger.tag(text6)
```

```
Out[45]: [('I', 'PPSS'),
          ('was', 'BEDZ'),
          ('happy', 'JJ'),
          (',', ','),
          ('but', 'CC'),
          ('so', 'RB'),
          ('was', 'BEDZ'),
          ('my', 'PP$'),
          ('enemy', 'NN'),
          ('.', '.')]

```

```
In [46]: # 8. c)
text7 = word_tokenize('So, how was the exam?')
bigram_tagger.tag(text7)
```

```
Out[46]: [('So', 'RB'),
          (',', ','),
          ('how', 'WRB'),
          ('was', 'BEDZ'),
          ('the', 'AT'),
          ('exam', None),
          ('?', None)]
```

```
In [47]: # 8. d)
text8 = word_tokenize('The students came in early so they can get good seats.')
bigram_tagger.tag(text8)
```

```
Out[47]: [('The', 'AT'),
          ('students', 'NNS'),
          ('came', 'VBD'),
          ('in', 'IN'),
          ('early', 'JJ'),
          ('so', 'CS'),
          ('they', 'PPSS'),
          ('can', 'MD'),
          ('get', 'VB'),
          ('good', 'JJ'),
          ('seats', 'NNS'),
          ('.', '.')] 
```

```
In [48]: # 8. e)
text9 = word_tokenize('She failed the exam, so she must take it again.')
bigram_tagger.tag(text9)
```

```
Out[48]: [('She', 'PPS'),
          ('failed', 'VBD'),
          ('the', 'AT'),
          ('exam', None),
          (',', None),
          ('so', None),
          ('she', None),
          ('must', None),
          ('take', None),
          ('it', None),
          ('again', None),
          ('.', None)]
```

```
In [49]: # 8. f)
text10 = word_tokenize('That was so incredible.')
bigram_tagger.tag(text10)
```

```
Out[49]: [('That', 'DT'),
          ('was', 'BEDZ'),
          ('so', 'QL'),
          ('incredible', 'JJ'),
          ('.', '.')]


```

```
In [50]: # 8. g)
text11 = word_tokenize('Wow, so incredible.')
bigram_tagger.tag(text11)
```

```
Out[50]: [('Wow', None), ('.', None), ('so', None), ('incredible', None), ('.', None)]


```

```
In [ ]:
```