

```
In [1]: 1 from nltk import word_tokenize
```

```
In [2]: 1 def preprocess(d):  
2     d=d.lower()  
3     d="eos "+ d  
4     d=d.replace(".", " eos")  
5     return d  
6 def generate_tokens(d):  
7     tokens = word_tokenize(d)  
8     return tokens  
9 def generate_tokens_freq(tokens):  
10    dct={}  
11    for i in tokens:  
12        dct[i]=0  
13    for i in tokens:  
14        dct[i]+=1  
15    return dct
```

In [3]:

```
1 def generate_ngrams(tokens,k):
2     l=[]
3     i=0
4     while(i<len(tokens)):
5         l.append(tokens[i:i+k])
6         i=i+1
7     l=l[:-1]
8     return l
9 def generate_ngram_freq(bigram):
10    dct1={}
11    for i in bigram:
12        st=" ".join(i)
13        dct1[st]=0
14    for i in bigram:
15        st=" ".join(i)
16        dct1[st]+=1
17    return dct1
```

In [4]:

```
1 def find1(s,dct1):
2     try:
3         return dct1[s]
4     except:
5         return 0
6 def print_probability_table(distinct_tokens,dct,dct1):
7     n=len(distinct_tokens)
8     l=[[]*n for i in range(n)]
9     for i in range(n):
10        denominator = dct[distinct_tokens[i]]
11        for j in range(n):
12            numerator = find1(distinct_tokens[i]+" "+distinct_tokens[j],dct1)
13            l[i].append(float("{:.3f}".format(numerator/denominator)))
14    return l
```

In []:

```
1 "The cat I saw sat on a mat was sad"
2 "The cat sat on a mat was sad"
```

Input

```
In [5]: 1 d=input("Enter corpus = ")
        2 print("\n"+'\033[1m'+ "Given Corpus" +'\033[0m')
        3 print(d)
```

Enter corpus = The cat sat on a mat was sad

Given Corpus

The cat sat on a mat was sad

Preprocess

```
In [6]: 1 d=preprocess(d)
        2 print("\n"+'\033[1m'+ "Preprocessing" +'\033[0m')
        3 print(d)
```

Preprocessing

eos the cat sat on a mat was sad

Generate Tokens

```
In [7]: 1 tokens=generate_tokens(d)
        2 print("\n"+'\033[1m'+ "Generate Tokens" +'\033[0m')
        3 print(tokens)
```

Generate Tokens

['eos', 'the', 'cat', 'sat', 'on', 'a', 'mat', 'was', 'sad']

Generate Tokens Frequency

```
In [8]: 1 distinct_tokens = list(set(sorted(tokens)))
        2 dct=generate_tokens_freq(tokens)
        3 print("\n"+'\033[1m'+ "Generate Frequency of Tokens" +'\033[0m')
        4 print(dct)
```

Generate Frequency of Tokens

```
{'eos': 1, 'the': 1, 'cat': 1, 'sat': 1, 'on': 1, 'a': 1, 'mat': 1, 'was': 1, 'sad': 1}
```

Generate bigram

```
In [9]: 1 bigram = generate_ngrams(tokens,2)
        2 print("\n"+'\033[1m'+ "Generate bigrams" +'\033[0m')
        3 for i in bigram:
        4     print("{} {}".format(' '.join(i)), end=", ")
```

Generate bigrams

```
'eos the', 'the cat', 'cat sat', 'sat on', 'on a', 'a mat', 'mat was', 'was sad',
```

Generate bigram frequency

```
In [10]: 1 dct1=generate_ngram_freq(bigram)
        2 print("\n\n"+'\033[1m'+ "Generate Frequency of bigrams" +'\033[0m')
        3 print(dct1)
```

Generate Frequency of bigrams

```
{'eos the': 1, 'the cat': 1, 'cat sat': 1, 'sat on': 1, 'on a': 1, 'a mat': 1, 'mat was': 1, 'was sad': 1}
```

Probability table

```
In [11]: 1 probability_table=print_probability_table(distinct_tokens,dct,dct1)
2 print("\n"+'\033[1m'+ "Probability table" +'\033[0m'+ "\n")
```

Probability table

```
In [12]: 1 n=len(distinct_tokens)
2 print("\t"+'\033[1m', end="")
3 for i in range(n):
4     print(distinct_tokens[i],end="\t")
5 print('\033[0m'+ "\n")
6
7 for i in range(n):
8     print('\033[1m',distinct_tokens[i],'\033[0m',end="\t")
9     for j in range(n):
10        print(probability_table[i][j],end="\t")
11    print("")
```

	a	mat	sad	the	was	cat	sat	on	eos
a	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
mat	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0
sad	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
the	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
was	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
cat	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0
sat	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
on	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
eos	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0

Testing

```
In [ ]: 1 "The cat sat on a mat"
        2 "The cat sad on a mat"
```

```
In [13]: 1 text = input("\nEnter text to check its probability = ")
        2 print("\n"+'\033[1m'+ "Given Text" +'\033[0m')
        3 print(text)
```

Enter text to check its probability = The cat sad on a mat

Given Text

The cat sad on a mat

```
In [14]: 1 p = preprocess(text)
        2 print("\n"+'\033[1m'+ "Preprocessing" +'\033[0m')
        3 print(p)
        4
        5 t=generate_tokens(p)
        6 print("\n"+'\033[1m'+ "Generate Tokens" +'\033[0m')
        7 print(t)
        8
        9 n = generate_ngrams(t,2)
       10 print("\n"+'\033[1m'+ "Generate bigrams" +'\033[0m')
       11 for i in n:
       12     print("{} {}".format(' '.join(i)), end=", ")
       13
```

Preprocessing

eos the cat sad on a mat

Generate Tokens

['eos', 'the', 'cat', 'sad', 'on', 'a', 'mat']

Generate bigrams

'eos the', 'the cat', 'cat sad', 'sad on', 'on a', 'a mat',

```

In [15]: 1 print("\n\n"+'\033[1m"+"Calculate bigram probability'+'\033[0m')
          2 s=1
          3 dct2={}
          4 for i in n:
          5     dct2[" ".join(i)]=0
          6
          7 for i in n:
          8     k=distinct_tokens.index(i[0])
          9     m=distinct_tokens.index(i[1])
         10
         11     dct2[" ".join(i)]=probability_table[k][m]
         12
         13     print("P('{})'\t= ".format(' '.join(i)),probability_table[k][m])
         14     s*=probability_table[k][m]

```

Calculate bigram probability

```

P('eos the')    =  1.0
P('the cat')     =  1.0
P('cat sad')     =  0.0
P('sad on')      =  0.0
P('on a')        =  1.0
P('a mat')       =  1.0

```

```

In [16]: 1 print("\n"+'\033[1m'+ "Calculate Probability of the sentence"+'\033[0m')
          2
          3 print(f"P('{text}')) \n= ",end="")
          4
          5 x=dct2.popitem()
          6
          7 for i in dct2:
          8     print(f"P('{i}'))", end=" * ")
          9
         10 print(f"P('{x[0]}')\n= ", end='')
         11
         12 for i in dct2:
         13     print(dct2[i], end=" * ")
         14
         15 print(x[1],"\n=",s)
         16
         17 print("\n"+'\033[1m'+f"Probability('{text}') = "+"{:.5f}".format(s))

```

Calculate Probability of the sentence

$P(\text{'The cat sad on a mat'})$

$= P(\text{'eos the'}) * P(\text{'the cat'}) * P(\text{'cat sad'}) * P(\text{'sad on'}) * P(\text{'on a'}) * P(\text{'a mat'})$

$= 1.0 * 1.0 * 0.0 * 0.0 * 1.0 * 1.0$

$= 0.0$

Probability('The cat sad on a mat') = 0.00000

In []:

1