

Assignment No. 1

Aim : To Learn to Calculate bigrams from the given Corpus and Calculate the probability Of Sentence.

Basic Knowledge:

- An N-gram language model predicts the probability of a given N-gram within any sequence of words in the language. If we have a good N-gram model, we can predict $p(w | h)$ – what is the probability of seeing the word w given a history of previous words h – where the history contains $n-1$ words.
- The **bigram model** approximates the probability of word given all the previous words, by using only the conditional probability of the last preceding word:

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1})$$

- The assumption that the probability of a word depends only on the previous word is called a *Markov assumption*.
- We have to train the machine to understand the meaning of the Sentence using the probabilistic Approach.
- It generally tries to understand the context of the given word according to the surrounding of the word(the word that comes after or before).
- We can calculate the Probability Using Chain Rule Using Probability.

$$P(w_1, w_2, \dots, w_n) = \prod (P(w_i | w_1, w_2, \dots, w_{i-1}))$$

- To calculate the probability of sentence from the above formula is not Practical so Markov Came with an Assumption.
- Markov Assumption:

$$P(w_1, w_2, \dots, w_n) = \prod (P(w_i | w_{i-k} \dots w_i))$$

- Probability is Found Using :

$$P(A/B) = P(A, B) / P(B)$$

1. Calculating Bigrams :

- A 2-gram (or bigram) is a two-word sequence of words, like “I love”, “love reading”, or “Analytics Vidhya”. And a 3-gram (or trigram) is a three-word sequence of words like “I love reading”, “about data science” or “on Analytics Vidhya”.
- Examples of Bigrams : This Book, My Book , Our Book , Your Book ,His Book , Her Book etc.

Example:

- Sentence : We need to book our tickets soon.

Bigrams From Above Sentence are :

We need , need to , to book , book our , our tickets ,
tickets soon .

- Sentence : They are going to play together.

Bigrams From Above Sentence are :

They are , are going , going to , to play , play together.

2. Calculating The Probability Of a Sentence:

- Estimating Bigram Probability:

$$P(w_i|w_{i-1}) = \text{Count}(w_{i-1}, w_i) / \text{Count}(w_{i-1})$$

Example:

<s> I am Sam <\s>

<s> Sam I am <\s>

<s> I do not like to study<\s>

Bigrams : <s>I, I am, am Sam , Sam <\s> , <s> Sam , Sam I , am <\s> , I do , do not, not like , like to , to study , study <\s>.

- Probabilities Of Bigrams :

$$P(I|<s>) = c(<s>,I)/c(<s>) = 2/3 = 0.67$$

$$P(\text{Sam}/<s>) = c(<s>,\text{sam})/c(<s>) = 1/3 = 0.33$$

$$P(\text{am}|I) = c(I,\text{am})/c(I) = 2/3 = 0.67$$

$$P(</s>|\text{Sam}) = 1/2 = 0.5$$

$$P(\text{Sam}|\text{am}) = 1/2 = 0.5$$

$$P(\text{do}|I) = 1/3 = 0.33.$$

```
In [1]: def ReadData():
in1 = input("Enter the Paragraph\n")
ls = in1.split('.')
ls1 = []
ls2 = []
for i in ls:
    if i == '':
        pass
    else:
        ls2.append(i)
for i in range(len(ls2)):
    if i == ' ':
        pass
    else:
        ls1.append(ls2[i].split())
return ls1

In [2]: def createBigram(data):
listOfBigrams = []
listOfUnigrams = []
bigramCounts = {}
unigramCounts = {}
for i in range(len(data)):
    for j in range(len(data[i])-1):
        listOfBigrams.append((data[i][j],data[i][j+1]))
        if bigramCounts.get((data[i][j],data[i][j+1])) == None:
            bigramCounts[(data[i][j],data[i][j+1])] = 1
        else:
            bigramCounts[(data[i][j],data[i][j+1])] = bigramCounts[(data[i][j],data[i][j+1])] + 1
for i in range(len(data)):
    for j in range(len(data[i])):
        listOfUnigrams.append(data[i][j])
        if unigramCounts.get(data[i][j]) == None:
```

```
            if unigramCounts.get(data[i][j]) == None:
                unigramCounts[data[i][j]] = 1
            else:
                unigramCounts[data[i][j]] = unigramCounts[data[i][j]] + 1
return listOfUnigrams,unigramCounts,listOfBigrams,bigramCounts
```

```
In [3]: def PrintUnigram(data):
listOfUnigrams = []
for i in range(len(data)):
    for j in range(len(data[i])):
        listOfUnigrams.append(data[i][j])
if listOfUnigrams == []:
    print("No Input Given")
    return
print(listOfUnigrams)
```

```
In [4]: def PrintBigram(data):
listOfBigrams = []
for i in range(len(data)):
    for j in range(len(data[i])-1):
        listOfBigrams.append((data[i][j],data[i][j+1]))
if listOfBigrams == []:
    print("No Input Given")
    return
print(listOfBigrams)
```

```
In [5]: def calcBigramProb(listOfBigrams,bigramCounts,unigramCounts):
listOfProb = {}
for bigram in listOfBigrams:
    word1 = bigram[0]
    word2 = bigram[1]
    listOfProb[bigram] = (bigramCounts.get(bigram)) / (unigramCounts.get(word1))
return listOfProb
```

```
In [6]: def Trained_model(listOfBigrams, bigramCounts, unigramCounts):
BigramProbability = calcBigramProb(listOfBigrams, bigramCounts, unigramCounts)
inlist = input("Write a Sentence to check the Probability\n")
inlist = inlist.split()
outputProb1 = 1
bilist = []
bigram = []
for i in range(len(inlist) - 1):
    if i < len(inlist) - 1:
        bilist.append((inlist[i], inlist[i + 1]))
print("The bigrams in given sentence are")
print(bilist)
for i in range(len(bilist)):
    if bilist[i] in BigramProbability:
        outputProb1 *= BigramProbability[bilist[i]]
    else:
        outputProb1 *= 0
print("The Probability Of the Sentence is ",outputProb1)
```

```
In [7]: def PrintAll(ls):
listOfUnigrams, unigramCounts, listOfBigrams, bigramCounts = createBigram(ls)
BigramProbability = calcBigramProb(listOfBigrams,bigramCounts,unigramCounts)
if listOfUnigrams == []:
    print("No Input Given")
    return
print("----list of Unigrams----")
print(listOfUnigrams)
print()
print("----list of Bigrams----")
print(listOfBigrams)
print()
print("----Frequency of Unigrams----")
print(unigramCounts)
```

```

print(unigramCounts)
print()
print("----Frequency of Bigrams----")
print(bigramCounts)
print()
print("----Probability of Bigrams----")
print(bigramProbability)
val = 1
ex = 0
lsl = [[]]
while(val==1 or ex <6):
    print('MENU : ')
    val = 0
    print(" 1: Writing the Paragraph \n 2: Print unigram \n 3: Print Bigram \n 4: Print Everything \n 5: Check The Trained Model")
    ex = int(input("Enter Your Choice: "))
    if ex == 1:
        lsl = ReadData()
        continue
    elif ex == 2:
        PrintUnigram(lsl)
        continue
    elif ex == 3:
        PrintBigram(lsl)
        continue
    elif ex == 4:
        PrintAll(lsl)
        continue
    elif ex == 5:
        listofUnigrams, unigramCounts, listofBigrams, bigramCounts = createBigram(lsl)
        Trained_model(listofBigrams, bigramCounts, unigramCounts)
        continue
    elif ex>6 or ex<=0:
        print("Program is Terminated")
        break

```

```

MENU :
1: Writing the Paragraph
2: Print unigram
3: Print Bigram
4: Print Everything
5: Check The Trained Model
6: Exit
Enter Your Choice: 1
Enter the Paragraph
Aakanksha Bhondve from comp B
MENU :
1: Writing the Paragraph
2: Print unigram
3: Print Bigram
4: Print Everything
5: Check The Trained Model
6: Exit
Enter Your Choice: 2
['Aakanksha', 'Bhondve', 'from', 'Comp', 'B']
MENU :
1: Writing the Paragraph
2: Print unigram
3: Print Bigram
4: Print Everything
5: Check The Trained Model
6: Exit
Enter Your Choice: 3
[['Aakanksha', 'Bhondve'], ('Bhondve', 'from'), ('from', 'Comp'), ('Comp', 'B')]
MENU :
1: Writing the Paragraph
2: Print unigram
3: Print Bigram
4: Print Everything
5: Check The Trained Model
6: Exit
Enter Your Choice: 4

```

```

Enter Your Choice: 4
----List Of Unigrams----
['Aakanksha', 'Bhondve', 'from', 'Comp', 'B']

----List Of Bigrams----
[['Aakanksha', 'Bhondve'], ('Bhondve', 'from'), ('from', 'Comp'), ('Comp', 'B')]

----Frequency Of Unigrams----
('Aakanksha': 1, 'Bhondve': 1, 'from': 1, 'Comp': 1, 'B': 1)

----Frequency of Bigrams----
(('Aakanksha', 'Bhondve'): 1, ('Bhondve', 'from'): 1, ('from', 'Comp'): 1, ('Comp', 'B'): 1)

----Probability of Bigrams----
(('Aakanksha', 'Bhondve'): 1.0, ('Bhondve', 'from'): 1.0, ('from', 'Comp'): 1.0, ('Comp', 'B'): 1.0)
MENU :
1: Writing the Paragraph
2: Print unigram
3: Print Bigram
4: Print Everything
5: Check The Trained Model
6: Exit
Enter Your Choice: 5
Write a Sentence to check the Probability
Comp B
The bigrams in given sentence are
[('Comp', 'B')]
The Probability Of the Sentence is 1.0
MENU :
1: Writing the Paragraph
2: Print unigram
3: Print Bigram
4: Print Everything
5: Check The Trained Model
6: Exit
Enter Your Choice: 6

```

```

0: exit
Enter Your Choice: 6
Program is Terminated

```