

# Assignment 02 - Language Modelling

Anirudh Kalla

October 14, 2020

## 1 Question 01

Vocabulary=100236

```
In [145]: print(vocabulary)
100236
```

Figure 1: Vocabulary

## 2 Question 02

Most Probable Unigrams:

1. jury
2. of

```
In [146]: unique_w = set(w)
largest_w = max(unique_w)
unique_w.remove(largest_w)
second_largest_w = max(unique_w)
print(largest_w)
print(second_largest_w)

0.061511025282282694
0.035368593217333844

In [147]: indices1 = [i for i,x in enumerate(w) if x == 0.061511025282282694]
indices2 = [i for i,x in enumerate(w) if x == 0.035368593217333844]
print(unl[indices1[0]])
print(unl[indices2[0]])

jury
of
```

Figure 2: Most Probable Unigrams

## 3 Question 03

The Most Frequent Bigram:-  
(‘of’, ‘the’)

```
In [148]: import statistics
from statistics import mode
mode(bi)

Out[148]: ('of', 'the')
```

Figure 3: Most Probable Bigrams

## 4 Question 04

Probability of an unseen Bigram :-  
 $9.976455564866915 \times 10^{-6}$

9.976455564866915e-06

Figure 4: P(Unseen Bigram)

## 5 Question 06.a

Probability of given sentence using the Vanilla MLE Bigram Model is:-  
0.0 This is because none of these bigrams occur in the brown-corpus.

Therefore according to Maximum Likelihood Estimation, this sentence has NULL probability.

```
Read 8 tokens
Generated 8 unigrams
{'<s>': 1, 'peter': 1, 'piper': 1, 'picked': 1, 'a': 1, 'pickled': 1, 'pepper': 1, '</s>': 1}
Generated 7 bigrams
{('<s>', 'peter'): 1, ('peter', 'piper'): 1, ('piper', 'picked'): 1, ('picked', 'a'): 1, ('a', 'pickled'): 1, ('pickled', 'pepper'): 1, ('pepper', '</s>'): 1}
0.0
```

Figure 5: P(Given Sentence: Vanilla MLE)

## 6 Question 06.b

[illegible]

This is because, even though none of these bigrams occur in the brown-corpus each of these UNK bigrams are assigned a (smoothed) probability. Giving the  $P(\text{Given Sentence})$  a non-zero value.

[illegible]

Figure 6: P(Given Sentence: Laplace smoothed)

## 7 Question 07

## MLE table of Bigrams (Vanilla)

```
x=[]
x0=[]
y=[]
data1=[]
data2=[]
data3=[]
data4=[]
for k,v in bigrams.items():
    first_word = k[0]
    first_word_count = unigrams[first_word]
    bi_prob = bigrams[k] / unigrams[first_word]
    bi_prob_laplace = (bigrams[k]+1) / (unigrams[first_word]+vocabulary)
    x.append(bi_prob)
    y.append(bi_prob_laplace)
    uni_prob = unigrams[first_word] / total_corpus
    w.append(uni_prob)
    final_prob = bi_prob * uni_prob
    print(k[0] + ' ' + x[1] + ' ' + y[1] + ' ' + str(v) + ' ' + 'lt' + str(first_word_count) + ' ' + 'lt' + 'lt' + str(bi_prob))
    data1.append(k[0] + ' ' + x[1])
    data2.append(str(v))
    data3.append(str(first_word_count))
    data4.append(str(bi_prob))
data = list(zip(data1, data2, data3, data4))
df = pd.DataFrame(data, columns = ['Bigrams', 'Bigrams Count', 'Unigram Count', 'Bigram Probability'])
df.head()
```

Read 1014940 tokens  
Generated 100216 unigrams  
Generated 539927 bigrams

	Bigrams	Bigrams Count	Unigram Count	Bigram Probability
0	The Fulton	1	6452	0.00015499070055796551
1	Fulton County	5	16	0.3125
2	County Grand	1	58	0.017241379310344827
3	Grand Jury	2	17	0.11764705882352941
4	Jury said	1	2	0.5

Figure 7: Bigram MLE Table

## 8 Question 08

MLE table of Bigrams (Laplace Smoothed)



Figure 8: Bigram MLE (Laplace Smoothed) Table