

# Assignment 1 Language Modelling

Prakash Kumar Uttam

SR -15247

<https://github.com/pkuttam/ngram-assignment>

## 1 Question - Perplexity Measure:

**Tr-G** = train on gutenber dataset (80%)

**Tr-B** = train on brown dataset (80%)

**Tr-GB** = train on gutenber and brown dataset (80%) each

**Ts-G** = test on gutenber dataset (20%)

**Ts-B** = test on brown dataset (20%)

more . give you a great multitude ,

**Example- Gutenberg dataset + brown dataset** i have not been able to go to the UNK of the lord , and the lord . been a great deal of the UNK , and he said , " i am sure i should have been the case of the house of the children of israel , and i

[here we can see UNK tag appear due to sparsity in brown dataset]

### 1.1 Add-K smoothing

—	uni-gram	bi-gram	trigram
<b>Tr-B and Ts-B</b>	399	291	Large
<b>Tr-G and Ts-G</b>	514	592	Large
<b>Tr-GB and Ts-B</b>	600	100	Large
<b>Tr-GB and Ts-G</b>	572	655	Large

### 1.2 kneser-Nay smoothing

—	uni-gram	bi-gram	trigram
<b>Tr-B and Ts-B</b>	-	63	97
<b>Tr-G and Ts-G</b>	-	98	160
<b>Tr-GB and Ts-B</b>	-	109	175
<b>Tr-GB and Ts-G</b>	-	75	110

## 2 Question - Sentence Generation :

### 2.1 add-K smoothing

**Example- brown dataset** i have been a UNK , and the UNK of the UNK , UNK , the UNK . to be a UNK of UNK , a UNK .

**Example- Gutenberg dataset**

i will not be afraid of the lord , and the lord . make thee a great deal of the house of the children of israel , and he said , " i am sure i should have been a great many

### 2.2 kneser-Nay smoothing

**Example- brown dataset**

i have been , and the UNK of the UNK , and UNK , UNK , the UNK . to be a UNK , a UNK .

**Example- Gutenberg dataset**

i will not be a great deal of the lord , and the lord . make thee a man of god , and he said , " i am sure i should be the lord god of israel , and i will give you a great , and to the lord

**Example- Gutenberg dataset + brown dataset**

i have been a UNK , and the lord , and he said , " i am sure i should be the lord . not been in the land of egypt , and i will not be a great deal of the lord god of israel , and to the lord

[here we can see UNK tag appear due to sparsity in brown dataset]

## 3 observation

The best model is Gutenberg + Brown data-set trained on tri-gram with Kneser-Nay smoothing.