# 1.0 Question 1

This question will be using the Text Corpus data to demonstrate the manual calculation and the calculation using Python programming via Colab platform on unsmoothed unigram probabilities and smoothed unigram probabilities.

## 1.1 Unsmoothed Unigram Probabilities

| Token | Frequency | Probability |
|-------|-----------|-------------|
| He | 2 | $P(He) = \dfrac{Count(He)}{Total\ Token}$ <br> $= \dfrac{2}{15}$ <br> $\approx 0.1333$ |
| read | 3 | $P(read) = \dfrac{Count(read)}{Total\ Token}$ <br> $= \dfrac{3}{15}$ <br> $= \dfrac{1}{5}$ <br> $= 0.2$ |
| a | 3 | $P(a) = \dfrac{Count(a)}{Total\ Token}$ <br> $= \dfrac{3}{15}$ <br> $= \dfrac{1}{5}$ <br> $= 0.2$ |
| book | 3 | $P(book) = \dfrac{Count(book)}{Total\ Token}$ <br> $= \dfrac{3}{15}$ <br> $= \dfrac{1}{5}$ <br> $= 0.2$ |
| I | 1 | $P(I) = \dfrac{Count(I)}{Total\ Token}$ <br> $= \dfrac{1}{15}$ <br> $\approx 0.0666$ |

| different | 1 | $P(\text{different}) = \dfrac{Count(different)}{Total\ Token}$ |
| --- | --- | --- |
| | | $= \dfrac{1}{15}$ |
| | | $\approx 0.0666$ |
| my | 1 | $P(\text{my}) = \dfrac{Count(my)}{Total\ Token}$ |
| | | $= \dfrac{1}{15}$ |
| | | $\approx 0.0666$ |
| Danielle | 1 | $P(\text{Danielle}) = \dfrac{Count(Danielle)}{Total\ Token}$ |
| | | $= \dfrac{1}{15}$ |
| | | $\approx 0.0666$ |

## 1.2 Smoothed Unigram Probabilities

| Token | Frequency | Probability |
| --- | --- | --- |
| He | 2 | $P(\text{He}) = \dfrac{Count(He)+1}{Total\ Token+Unique\ Token+UNK}$ |
| | | $= \dfrac{2+1}{15+8+1}$ |
| | | $= \dfrac{3}{24}$ |
| | | $= \dfrac{1}{8}$ |
| | | $= 0.125$ |
| read | 3 | $P(\text{read}) = \dfrac{Count(read)+1}{Total\ Token+Unique\ Token+UNK}$ |
| | | $= \dfrac{3+1}{15+8+1}$ |
| | | $= \dfrac{4}{24}$ |
| | | $= \dfrac{1}{6}$ |
| | | $\approx 0.1666$ |
| a | 3 | $P(\text{a}) = \dfrac{Count(a)+1}{Total\ Token+Unique\ Token+UNK}$ |
| | | $= \dfrac{3+1}{15+8+1}$ |
| | | $= \dfrac{4}{24}$ |

| | | |
|---|---|---|
| | | $$= \frac{1}{6}$$ $$\approx 0.1666$$ |
| book | 3 | $$P(\text{book}) = \frac{Count(\text{book})+1}{Total\ Token+Unique\ Token+UNK}$$ $$= \frac{3+1}{15+8+1}$$ $$= \frac{4}{24}$$ $$= \frac{1}{6}$$ $$\approx 0.1666$$ |
| I | 1 | $$P(\text{I}) = \frac{Count(I)+1}{Total\ Token+Unique\ Token+UNK}$$ $$= \frac{1+1}{15+8+1}$$ $$= \frac{2}{24}$$ $$= \frac{1}{12}$$ $$\approx 0.0833$$ |
| different | 1 | $$P(\text{different}) = \frac{Count(different)+1}{Total\ Token+Unique\ Token+UNK}$$ $$= \frac{1+1}{15+8+1}$$ $$= \frac{2}{24}$$ $$= \frac{1}{12}$$ $$\approx 0.0833$$ |
| my | 1 | $$P(\text{my}) = \frac{Count(my)+1}{Total\ Token+Unique\ Token+UNK}$$ $$= \frac{1+1}{15+8+1}$$ $$= \frac{2}{24}$$ $$= \frac{1}{12}$$ $$\approx 0.0833$$ |
| Danielle | 1 | $$P(\text{Danielle}) = \frac{Count(Danielle)+1}{Total\ Token+Unique\ Token+UNK}$$ $$= \frac{1+1}{15+8+1}$$ $$= \frac{2}{24}$$ |

| | | $$= \frac{1}{12}$$ $$\approx 0.0833$$ |
| --- | --- | --- |

# 1.3 Python Implementation

## Reading Data

```
# Read data in text file
file = open("Text Corpus.txt")
data = file.read()
print(data)

<s> He read a book </s>
<s> I read a different book </s>
<s> He read a book my Danielle </s>
```

## Data Cleaning

```
# Tokenize every word in the data and clean the data
tokens = data.split()

cleanedTokens = [token for token in tokens if token != "<s>"]
cleanedTokens = [token for token in cleanedTokens if token != "</s>"]
print(cleanedTokens)

totalToken = len(cleanedTokens)

['He', 'read', 'a', 'book', 'I', 'read', 'a', 'different', 'book', 'He', 'read', 'a', 'book', 'my', 'Danielle']
```

## Unique Tokens

```
# Create a list of unique token from the data
uniqueTokens = list(set(cleanedTokens))
print(uniqueTokens)

totalUniqueToken = len(uniqueTokens)

['my', 'Danielle', 'read', 'different', 'I', 'book', 'He', 'a']
```

## Calculating Frequency

```
# Frequency for each unique token
frequencyResult = ""

for uniqueToken in uniqueTokens:
  frequency = cleanedTokens.count(uniqueToken)
  frequencyResult = frequencyResult + "Count(" + uniqueToken + ") = " + str(frequency) + "\n"
```

## Calculating Unsmoothed Unigram Probabilities

```python
# Unsmoothed Unigram Probabilities
unsmoothedResult = ""

for uniqueToken in uniqueTokens:
    probability = cleannedTokens.count(uniqueToken) / totalToken
    unsmoothedResult = unsmoothedResult + "P(" + uniqueToken + ") = " + str(probability) + "\n"
```

## Calculating Smoothed Unigram Probabilities

```python
# Smoothed Unigram Probabilities
smoothedResult = ""

for uniqueToken in uniqueTokens:
    probability = (cleannedTokens.count(uniqueToken) + 1) / (totalToken + totalUniqueToken + 1)
    smoothedResult = smoothedResult + "P(" + uniqueToken + ") = " + str(probability) + "\n"
```

## Output for Summary of Unigram Probabilities

```python
# Summary of unigram probabilities
print("Number of Tokens: " + str(totalToken))
print("Number of Unique Tokens: " + str(totalUniqueToken) + "\n")
print("Frequency of Tokens:\n" + frequencyResult)
print("Unsmoothed Probabilities of Tokens:\n" + unsmoothedResult)
print("Smoothed Probabilities of Tokens:\n" + smoothedResult)
```

```
Number of Tokens: 15
Number of Unique Tokens: 8

Frequency of Tokens:
Count(my) = 1
Count(Danielle) = 1
Count(read) = 3
Count(different) = 1
Count(I) = 1
Count(book) = 3
Count(He) = 2
Count(a) = 3

Unsmoothed Probabilities of Tokens:
P(my) = 0.06666666666666667
P(Danielle) = 0.06666666666666667
P(read) = 0.2
P(different) = 0.06666666666666667
P(I) = 0.06666666666666667
P(book) = 0.2
P(He) = 0.13333333333333333
P(a) = 0.2

Smoothed Probabilities of Tokens:
P(my) = 0.08333333333333333
P(Danielle) = 0.08333333333333333
P(read) = 0.16666666666666666
P(different) = 0.08333333333333333
P(I) = 0.08333333333333333
P(book) = 0.16666666666666666
P(He) = 0.125
P(a) = 0.16666666666666666
```