



INSTITUTE OF TECHNOLOGY OF CAMBODIA

Department of Information and Communication Engineering

MINI PROJECT 1 REPORT

Group 07

TOPIC: TEXT GENERATION

Deadline 21 December 2023

: Natural Language Processing Course

: VALY Dona Lecturer

Team Member

Name	ID
SOR Sopheak	e20191066
SUN Vysing	e20191124
THUE Naratt	e20191186
TY Sophearum	e20191219
VEN Thon	e20191250

I. Pre-Processing

1. Split paragraphs in the dataset into arrays of sentences.

```
# Tokenize sentences
from nltk.tokenize import sent_tokenize, word_tokenize
sentences = sent_tokenize(text)
print(sentences)
['The weather was cold and gray as usual at this time of year.', 'The trees were all leafless, with fall now just a memory.'
```

2. Clean data in each sentence to follow some set of rules and split words in a sentence into arrays of words.

```
# Function to tokenize words in sentences

def tokenize_word(sentences):
    sentences_tokenized = []
    for sentence in sentences:
        words = sentence.split()
        sentences_tokenized.append(words)
    return sentences_tokenized

# Tokenize words in sentences
sentences = tokenize_word(sentences)

['<s>', 'my', 'dream', 'house', 'would', 'be', 'located', 'in', 'the', 'country', ',', 'but', 'not', 'too', 'far', 'from', '
[<s>', 'since', 'the', 'basement', 'is', 'fairly', 'shallow', ',', 'there', 'are', 'three', 'foot', 'by', 'five', 'foot', '
['<s>', 'inow', ', 'imagine', 'that', 'another', 'person', 'who', 'was', 'also', 'at', 'the', 'store', 'saw', 'the', 'crim
['<ss', 'feeling', 'secure', 'means', 'being', 'at', 'home', ',', 'no', 'matter', 'the', 'atmosphere', ',', with', 'the', '
['<ss', 'right', 'now', ',', 'the', 'world', 'outside', 'of', 'home', 'was', 'more', 'safe', 'and', 'structured', ',', 'not'
['<s>', 'as', 'if', 'father', 'christmas', 'had', 'came', 'early', ',', 'their', 'faces', 'lit', 'up', 'as', 'if', 'they', 'ss', 'surly', 'the', 'london', 'eye', 'can't'', 'be', 'bigger', 'than', 'this', 'right?', '</s>']
['<ss', 'surly', 'the', 'infith', 'lasted', 'for', 'about', '5', 'minutes', 'uniti', 'a', 'police', 'man', 'came', 'and', 'arrested', 'cs', 'the', 'fight', 'lasted', 'for', 'about', '5', 'minutes', 'uniti', 'a', 'police', 'man', 'came, 'and', 'arrested', 'css', 'the', 'fight', 'lasted', 'for', 'about', '5', 'minutes', 'uniti', 'a', 'police', 'man', 'came, 'and', 'arrested', 'css', 'kids', 'were', 'crowding', 'around', 'so', 'they', 'could', 'get', 'first', 'pick', 'of', 'the', 'sugary', 'goodnes
['<ss', 'kids', 'were', 'crowding', 'around', 'so', 'they', 'could', 'get', 'first', 'pick', 'of', 'the', 'sugary', 'goodnes
['<ss', 'the', 'boomig', 'out', 'like', 'an', 'elephant', 'running', 'through', 'a', 'jungle', 'and', 'then', '"sos', 'less', 'wou', 'answers', 'to', 'these', 'questions', 'will', 'form', 'a', 'basis', 'for', 'the', 'qualities', 'and', 'char
['<ss', 'the', 'weathe
```

3. Split the corpus into 3 for train valid and test by 70% 10% and 20% respectively.

```
# Split data into 3
train_size = int(len(sentences)*0.7)
valid_size = int(len(sentences)*0.1)

train = sentences[:train_size]
valid = sentences[train_size:train_size+valid_size]
test = sentences[train_size+valid_size:]
print(len(train), len(valid), len(test))
237 34 69
```

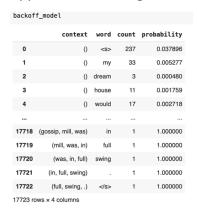
II. Build models

1. Back-off model

I build the back-off by doing the following steps:

a) Create a function to create any n-gram model by splitting the i to n-1 word in the sentence into the context column and taking the nth word to the word column and i

- start from 0 to end of sentence. Then it count the amount of time the unique context and word appear and calculating the probability.
- b) Use the function in a) to create 4 n-gram model (uni-gram, bi-gram, tri-gram and four-gram)
- c) Append all those n-gram models together and call it backoff_model:



2. Interpolation model

I build the interpolation by doing the following steps:

a) Create a function to create any n-gram model by splitting the i to n-1 word in the sentence into the context column and taking the nth word to the word column and i start from 0 to end of sentence. Then it count the amount of time the unique context and word appear and calculating the probability. But this time, use add-k smoothing to get the probability.

$$P_{Add-k}(w_i|w_{i-1}) = \frac{C(w_{i-1},w_i) + k}{C(w_{i-1}) + kV}$$

Figure 1: Add-k smoothing formula

- b) Use the function in a) to create 4 n-gram model (uni-gram-smooth, bi-gram-smooth, trigram-smooth and four-gram-smooth)
- c) In each model, add the probability of all its lower grams together and multiply each probability with its lamda values.
- d) Append all those n-gram models together and call it interpolation model:

interpolation_model				
	context	word	count	probability
0	0	<s></s>	237	0.029314
1	0	my	33	0.004188
2	0	dream	3	0.000493
3	0	house	11	0.001478
4	0	would	17	0.002217
17718	(gossip, mill, was)	in	1	0.003116
17719	(mill, was, in)	full	1	0.000978
17720	(was, in, full)	swing	1	0.000885
17721	(in, full, swing)		1	0.008341
17722	(full, swing, .)		1	0.036960
17723 rows × 4 columns				

III. Calculate perplexity

```
def perplexity(model, data, n):
    perplex = Decimal(1)
    loop_count = 0
    for sentence in data:
        for i in range(len(sentence)-1-n):
            loop_count = loop_count + 1
            perplex = perplex * ( Decimal(1) / Decimal(get_proba(model, tuple(sentence[i:n+i-1]), sentence[n+i-1])) )
        return perplex
    return perplex
```

Figure 2: Perplexity calculation function

By writing the above function, following the perplexity formula, we are able to calculate the perplexity of both model:

```
perplexity(backoff_model, test, 4)

Decimal('276.0602328431615391442621491')

Decimal('635.3798439510164544222046843')
```

Figure 3: interpolation model perplexity

Figure 4: back-off model perplexity

IV. Text generator

```
def generate(model, context=[]):
    sentence = ["<s>"]+context
    while sentence[-1] != "</s>" and len(sentence)<100:
        proba = get_proba_distrib(model, sentence)

    if sum(proba['probability'].tolist()) < 1:
        normalized_probabilities = proba['probability'].tolist() / np.sum(proba['probability'].tolist())

    w = np.random.choice(proba['word'].tolist(), 1, p = normalized_probabilities)
    sentence.append(w[0])

result = ''
for word in sentence:
    result = result + ' ' + word
return result</pre>
```

Figure 5: Text generator function

Write a function call generate, it find the next possible with the most probability to appear, we can create a text generator with both models.