# Assignment #5

CSE 447: Natural Language Processing
Eric Boris: 1976637
Collaborators: Allyson Ely, Ardi Madadi

## Implementation Problem

**Problem 2: Answer**

1. **Yes**, because if there are multiple words per line then different words will be joined rather than their letters being joined.

2. **No**, each additional iteration over the data produces a new BPE encoding.
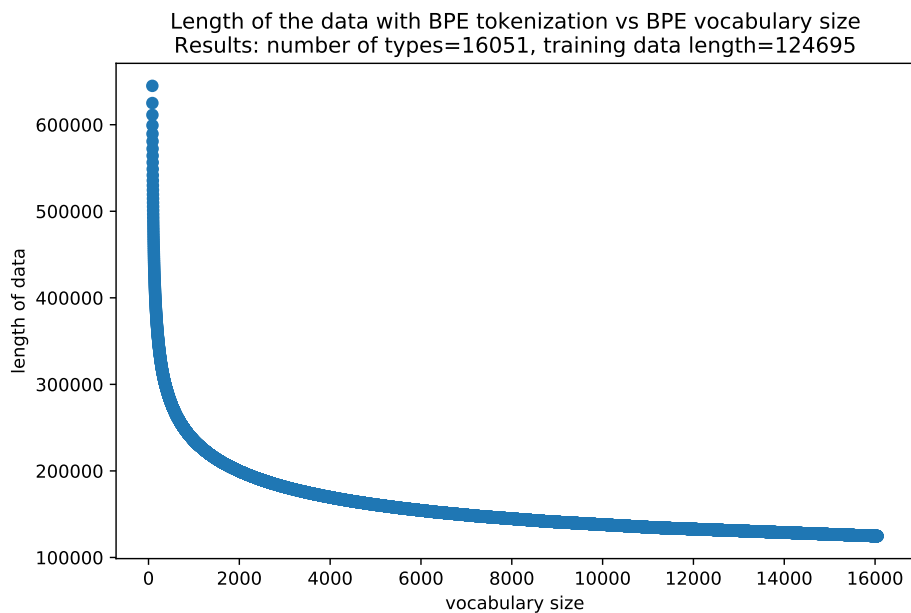
3. See Figure 1



Figure 1

4. **No**, we don't get any <unk> tokens. We're unlikely to get <unk> because that will only happen if the entire word cannot be encoded by BPE which will only occur if the sequence of letters composing that word hasn't been encoded during training. Using large, diverse training datasets reduces the likelihood of this occuring.

Results of encoding 10 uncommon english words using BPE

|    | Word | Result |
|----|------|--------|
| 1  | fudgel | f ud g el\<s\> |
| 2  | nudiustertian | nu di ust er tian\<s\> |
| 3  | selcouth | sel cou th\<s\> |
| 4  | zabernism | z ab ern ism\<s\> |
| 5  | douceur | dou ce ur\<s\> |
| 6  | pauciloquent | pa uc il o qu ent\<s\> |
| 7  | defenestration | defen estr ation\<s\> |
| 8  | hiraeth | h ir ae th\<s\> |
| 9  | limerence | li mer ence\<s\> |
| 10 | sonder | s ond er\<s\> |

5. Compared to this implementation of BPE, character encoding is quite fast since there are fewer comparisons. But it won't perform as well a BPE encoding because it's too fine grained and can't recognize word structures.

## Problem 2: Code

```
1   # CSE 447 A5 Problem 2
2
3   from collections import Counter
4   import matplotlib.pyplot as plt
5   import numpy as np
6
7   def load(path):
8       '''
9       Return the contents of the given file located in path.
10      '''
11      with open(path, 'r') as f:
12          return f.read()
13
14  def save(path, contents):
15      '''
16      Save the contents to the file at the given path.
17      '''
18      with open(path, 'w') as f:
19          f.write(contents)
20
21  def plot(title, subtitle, x_label, y_label, x, y, path, dim=(8, 5)):
22      '''
23      Plot the given data as a scatter plot.
24      '''
25      plt.figure(figsize=dim)
26      plt.title(f'{title}\n{subtitle}')
27      plt.xlabel(x_label)
28      plt.ylabel(y_label)
29      plt.plot(x, y, 'o')
30      plt.savefig(path)
31      plt.show()
32
33  class BPE:
34      def __init__(self):
35          self.end_sym = '<s>'
36
37      def get_vocab(self, data):
38          '''
39          Return the unique characters appearing in the data plus the end symbol.
40          '''
41          vocab = set(data)
42          vocab.remove(' ')
43          vocab = list(vocab)
44          vocab.append(self.end_sym)
45          return vocab
46
47      def tokenize(self, data):
```

```
48                ,,,
49            Return the data as a tokenized list of characters with words separated
50            by the end symbol.
51                ,,,
52            tokens = []
53            for word in data.split():
54                for char in word:
55                    tokens.append(char)
56                tokens.append(self.end_sym)
57            return tokens
58
59        def get_pairs(self, tokens):
60                ,,,
61            Return a mapping of token pairs -> count of token pair occurences.
62                ,,,
63            pairs = Counter()
64            for i in range(len(tokens)-1):
65                # Prevent merge across words.
66                if not tokens[i].endswith(self.end_sym):
67                    pairs[tokens[i], tokens[i+1]] += 1
68            return pairs
69
70        def merge(self, tokens, best):
71                ,,,
72            Merge all the best token pairs in tokens and return the merged list.
73                ,,,
74            i = 0
75            while i < len(tokens) - 1:
76                if (tokens[i], tokens[i+1]) == best:
77                    tokens[i : i+2] = [''.join(tokens[i : i+2])]
78                i += 1
79            return tokens
80
81        def train(self, data, min_freq=1, verbose=False, print_every=10):
82                ,,,
83            Train the bpe model and return the resultant vocabulary and merged
84            tokens and the size of the vocab and tokens at each step.
85                ,,,
86            self.vocab = self.get_vocab(data)
87            tokens = self.tokenize(data)
88            pairs = self.get_pairs(tokens)
89
90            # Let this be the tuple from pairs with the highest occurence.
91            best = max(pairs, key=pairs.get)
92
93            # Store the best transformations for applying BPE to new words
94            # after training.
95            self.transforms = []
96
97            # These will be returned and used for plotting data.
98            n_vocab = [len(self.vocab)]
99            n_tokens = [len(tokens)]
100
101            # Let this maintain the current iteration for displaying progress.
102            i = 0
103
104            while pairs[best] > min_freq:
105                self.vocab.append(''.join(best))
106                self.transforms.append(best)
107
108                tokens = self.merge(tokens, best)
109                pairs = self.get_pairs(tokens)
110                best = max(pairs, key=pairs.get)
111
112                n_vocab.append(len(self.vocab))
113                n_tokens.append(len(tokens))
114
115                if verbose and i % print_every == 0:
116                    print(f'{i} {pairs[best]}')
```

```python
117
118                 i += 1
119
120             return (' ').join(tokens), self.vocab, n_tokens, n_vocab
121
122     def encode(self, data):
123         '''
124         Apply the tranformations found during training to encode the given data
125         into the trained BPE scheme.
126         '''
127         tokens = self.tokenize(data)
128         for best in self.transforms:
129             tokens = self.merge(tokens, best)
130         return (' ').join(tokens)
131
132 if __name__ == '__main__':
133     # Load the training data.
134     load_path = '../data/A5-data.txt'
135     data = load(load_path)
136
137     # Build and train the model.
138     bpe = BPE()
139     encoded, vocab, n_tokens, n_vocab = bpe.train(data, verbose=True, print_every=100)
140
141     # Display the results.
142     print(f'encoded={encoded}\nvocab={vocab}\nn_tokens={n_tokens}\nn_vocab={n_vocab}')
143
144     # Save the encoded output.
145     save_path = '../data/output.txt'
146     save(save_path, encoded)
147
148     # Plot the effects of training on data length and vocab size.
149     plot(title='Length of the data with BPE tokenization vs BPE vocabulary size',
150         subtitle=f'Results: number of types={n_vocab[-1]}, training data length={n_tokens[-1]}
             ',
151         x_label='vocabulary size',
152         y_label='length of data',
153         x=n_vocab,
154         y=n_tokens,
155         path='../figures/plot.pdf')
156
157     # Encode new words using the trained BPE model.
158     new_load_path = '../data/least_common_words.txt'
159     new_words = load(new_load_path)
160     new_encoded = bpe.encode(new_words)
161
162     # Display the results.
163     print(new_encoded)
164
165     # Save the new encoded output.
166     new_save_path = '../data/new_output.txt'
167     save(new_save_path, new_encoded)
```