

```
In [103]: 1 import numpy as np
          2 import pandas as pd
          3 import argparse
          4 import math
          5 import random
          6 from tqdm import tqdm
          7 import matplotlib.pyplot as plt
          8 import os
          9 import operator
         10 from collections import defaultdict as dd
         11
```

```
In [101]: 1 masked_txt = pd.read_csv (r'/Users/nehakardam/Documents/UWclasses /CSE
          2 masked_txt.to_csv (r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A6/
          3 lm_txt = pd.read_csv (r'/Users/nehakardam/Documents/UWclasses /CSE NLP/
          4 lm_txt.to_csv (r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A6/data
```

<ipython-input-101-32cd5041a44b>:1: ParserWarning: Falling back to the 'python' engine because the 'c' engine does not support regex separators (separators > 1 char and different from '\s+' are interpreted as regex); you can avoid this warning by specifying engine='python'.

```
masked_txt = pd.read_csv (r'/Users/nehakardam/Documents/UWclasses /CSE
NLP/A6/data_A6/15pctmasked.txt', sep='delimiter', names = ["Sentences"])
```

<ipython-input-101-32cd5041a44b>:3: ParserWarning: Falling back to the 'python' engine because the 'c' engine does not support regex separators (separators > 1 char and different from '\s+' are interpreted as regex); you can avoid this warning by specifying engine='python'.

```
lm_txt = pd.read_csv (r'/Users/nehakardam/Documents/UWclasses /CSE NLP/
A6/data_A6/lm.txt', sep='delimiter', names = ["Words", "Prob"])
```

```
In [109]: 1 blm = BigramModel('/Users/nehakardam/Documents/UWclasses /CSE NLP/A6/da
          2 v = Viterbi('/Users/nehakardam/Documents/UWclasses /CSE NLP/A6/data_A6/
          3                '/Users/nehakardam/Documents/UWclasses /CSE NLP/A6/data_A6/
          4 complete_sentences = v.compute_missing_characters(blm)
          5 v.write_sentences_to_file(complete_sentences)
```

Running Viterbi Algorithm: 100%|██████████| 5785/5785 [04:58<00:00, 19.36 it/s]

```
In [129]: 1 print(blm)
```

<__main__.BigramModel object at 0x7f87c8fbbf40>

In [116]:

```

1 START = '<start>'
2 MASK = '<mask>'
3 SPACE = '<s>'
4 EOS = '<eos>'
5
6 class BigramModel:
7     def __init__(self, file_path):
8         if not os.path.exists(file_path):
9             print(f"Path: {file_path} does not exist")
10            return None
11
12        with open(file_path, "r") as f:
13            bigrams = f.readlines()
14
15        self.blm = defaultdict(dict)
16
17        for bigram in bigrams:
18            bigram, prob = bigram.split('\t')[0], bigram.split('\t')[1]
19
20            w1, w2 = bigram.split(' ')[0], bigram.split(' ')[1]
21            prob = float(prob.strip())
22
23            self.blm[w1][w2] = prob
24
25        self._create_labels_dict()
26
27    def _create_labels_dict(self):
28        labels = list(self.blm.keys())
29        labels.sort()
30        labels.append(labels.pop(labels.index(EOS))) # add EOS to end
31        labels = [labels.pop(labels.index(START))] + labels # add START
32
33        self.labels = labels
34        self.label_to_idx = {k: v for v, k in enumerate(labels)}
35        self.idx_to_label = {v: k for v, k in enumerate(labels)}
36
37    def get_labels(self):
38        return self.labels
39
40    def get_labels_to_index(self):
41        return self.label_to_idx
42
43    def get_index_to_labels(self):
44        return self.idx_to_label
45
46    def get_bigram_prob(self, w1, w2) -> float:
47        if w1 in self.blm.keys() and w2 in self.blm[w1].keys():
48            return self.blm[w1][w2] # p(w2 | w1) = prob
49        return 0
50
51    def get_w2_given_w1(self, w1, w2, is_log_prob=True):
52        if w1 in self.blm.keys() and w2 in self.blm[w1].keys():
53            if is_log_prob: return np.log(self.blm[w1][w2]) # p(w2 | w1) = ln(prob)
54            else: return self.blm[w1][w2] # p(w2 | w1) = prob
55
56        if is_log_prob: return float('-inf')

```

```

57         return 0
58
59     def get_max_from_key(self, w1) -> str:
60         return max(self.blm[w1].items(), key=operator.itemgetter(1))
61
62 class Viterbi:
63     def __init__(self, input_file_path, output_file_path):
64         if not os.path.exists(input_file_path):
65             print(f"Path: {input_file_path} does not exist")
66             return None
67
68         # Read and parse the input file
69         with open(input_file_path, "r") as f:
70             l = f.readlines()
71
72         self.masked_sentences = []
73         for sentence in l:
74             sentence = sentence.split()
75             self.masked_sentences.append(sentence)
76
77         self.output_file_path = output_file_path
78
79     def compute_missing_characters(self, blm: BigramModel):
80         states = blm.get_labels()
81         complete_sentences = []
82
83         for sentence in tqdm(self.masked_sentences, desc="Running Vite:
84             best_path = self.viterbi_algorithm(sentence, states, blm)
85             complete_sentences.append(best_path)
86
87         return complete_sentences
88
89     def viterbi_algorithm(self, observation, states, blm: BigramModel)
90
91         sentence = observation[:]
92
93         idx_to_state = blm.get_index_to_labels()
94         state_to_idx = blm.get_labels_to_index()
95         R, C = len(states), len(sentence)
96
97         # To hold p. of each state given each sentence.
98         trellis = np.full((R, C), -np.inf)
99
100        # to hold the back pointers for cell
101        back_pointer = np.zeros((R, C), dtype='int32')
102
103        # Determine each hidden state's p. at time 0
104        for i in range(R-1):
105            if states[i] == START:
106                trellis[i][0] = blm.get_w2_given_w1(START, START) # in
107
108        # and now, assuming each state's most likely prior state, k
109        for j in range(1, C-1):
110            w1, w2 = sentence[j-1], sentence[j]
111
112            for i in range(R-1):
113                label = states[i]

```

```

114
115     # Case 1: w1 and w2 are both known characters
116     if w1 != MASK and w2 != MASK:
117         label_idx = state_to_idx[w2]
118         max_prev_trellis_value = max(trellis[:, j-1])
119         max_prev_trellis_label_idx = np.argmax(trellis[:, j-1])
120
121         trellis[label_idx][j] = max_prev_trellis_value + blm.get_w1(w1)
122         back_pointer[label_idx][j] = max_prev_trellis_label_idx
123         break
124
125     # Case 2: curr is MASK and prev column is known
126     elif w1 != MASK and w2 == MASK:
127
128         max_prev_trellis_value = max(trellis[:, j-1])
129         max_prev_trellis_label_idx = np.argmax(trellis[:, j-1])
130
131         trellis[i][j] = max_prev_trellis_value + blm.get_w2(w2)
132         back_pointer[i][j] = max_prev_trellis_label_idx
133
134     # Case 3: curr is known and prev column is MASK
135     elif w2 != MASK and w1 == MASK:
136         label_idx = state_to_idx[w2]
137
138         t1 = np.full((R-1, ), -np.inf)
139         for k in range(R-1):
140             prev_label = states[k]
141             t1[k] = blm.get_w2_given_w1(w1=prev_label, w2=w2)
142
143         trellis[label_idx][j] = max(t1)
144         back_pointer[label_idx][j] = np.argmax(t1)
145         break
146
147     # Case 4: w1 and w2 are both MASK characters
148     else:
149         t1 = np.full((R-1, ), -np.inf)
150
151         '''
152         Since the prev column is also a mask, we pick the max from
153         all the other labels as well. Then find the max
154         '''
155         for k in range(R-1):
156             prev_label = states[k]
157             t1[k] = blm.get_w2_given_w1(w1=prev_label, w2=w2)
158         trellis[i][j] = max(t1)
159         back_pointer[i][j] = np.argmax(t1)
160
161     # Fill in the prob for <eos>
162     trellis[R-1][C-1] = blm.get_w2_given_w1(w1=EOS, w2=EOS)
163     back_pointer[R-1][C-1] = np.argmax(trellis[:, C-2])
164
165     # np.savetxt('back_pointer.out', np.vstack(['header'] + sentence, 'back_pointer.out'))
166     # np.savetxt('trellis.out', np.vstack(['header'] + sentence, 'trellis.out'))
167
168     # get the back pointers
169     guessed_sentence = [EOS]
170     label_idx = back_pointer[R-1][C-1]

```

```

171         for j in reversed(range(C-1)):
172             guessed_sentence.append(states[label_idx])
173             label_idx = back_pointer[label_idx][j]
174
175         return list(reversed(guessed_sentence))
176
177     def write_sentences_to_file(self, sentence_list):
178         sentence_list_strings = []
179
180         for sentence in sentence_list:
181             s = ' '.join(sentence)
182             sentence_list_strings.append(s)
183
184         s = '\n'.join(sentence_list_strings)
185
186         with open(self.output_file_path, "w") as f:
187             f.write(s)
188
189     def sanity_check_output(masked_sentences, un_masked_sentences):
190         print("Performing sanity check on output")
191
192         # Ensure the first sentence matches the correct output:
193         correct_out = ['<start>', 'I', '<s>', 'p', 'e', '<s>', 'm', 'a', 't']
194         for correct_char, unmasked_char in zip(correct_out, un_masked_sentences):
195             if correct_char != unmasked_char:
196                 print(f'Error! Viterbi Algorithm is incorrect for sentence')
197                 return
198
199         for masked_sentence, unmasked_sentence in zip(masked_sentences, un_masked_sentences):
200
201             # Ensure the length of 2 sentence is the same (same number of tokens)
202             lm, lum = len(masked_sentence), len(unmasked_sentence)
203             if lm != lum:
204                 print(f'Error! The length of the sentences do not match')
205                 print(f'Masked Sentence: {masked_sentence}')
206                 print(f'Unmasked Sentence: {unmasked_sentence}')
207
208             # Ensure we only changed the <mask> characters
209             for i in range(lm):
210                 c_m, c_um = masked_sentence[i], unmasked_sentence[i]
211
212                 if c_m != MASK and (c_m != c_um):
213                     print(f'Error! Changed a known character')
214                     print(f'Changed {c_m} -> {c_um} at index: {i}')
215
216                 elif c_m == MASK and c_um == START:
217                     print(f'Changed a masked char to <start> token!')
218                     print(masked_sentence)
219                     print(unmasked_sentence)
220                     print("")
221
222     def parse_output_file(output_file_path):
223         un_masked_sentences = []
224
225         # Read in output file if it exists
226         if os.path.exists(output_file_path):
227             with open(output_file_path, "r") as f:

```

```
228         l = f.readlines()
229         for sentence in l:
230             sentence = sentence.split()
231             un_masked_sentences.append(sentence)
232
233         return un_masked_sentences
234
235 # if __name__ == "__main__":
236 #     print("NLP - A4")
237
238 #     parser = argparse.ArgumentParser(description='Viterbi Algorithm')
239 #     parser.add_argument("-lm", "--lang-model", dest="lang_model_path")
240 #     parser.add_argument("-ip", "--input-file", dest="input_file_path")
241 #     parser.add_argument("-op", "--output-file", dest="output_file_path")
242 #     parser.add_argument("-t", "--sanity-check", dest="perform_sanity_check")
243 #     args = parser.parse_args()
244
245 #     blm = BigramModel(args.lang_model_path)
246 #     v = Viterbi(input_file_path=args.input_file_path, output_file_path=args.output_file_path)
247
248 #     complete_sentences = v.compute_missing_characters(blm)
249 #     v.write_sentences_to_file(complete_sentences)
250
251 #     if args.perform_sanity_check:
252 #         sanity_check_output(v.masked_sentences, parse_output_file(args.output_file_path))
```

In []: 1