

# Trigram Model: Theory and Practice

## 1 Preprocessing each line

Our main task is to train a character-based language model identify different languages. Inconsistent symbols, or redundant punctuations do not convey much information and will bias our identification task; therefore, we intend to preprocess each training text by lines, retaining only English Alphabets, spaces, digits (all converted to 0) and periods. In addition, we reckon extra attention paid to the new line character ‘\n’, which can be easily mistaken for a space.

Based on phonotactics and morphotactics of natural languages, some morphemes or phoneme combinations are inhibited; naturally, some characters are not likely to occur at the beginning or the end of a sentence. To capture this property, we add start symbols and end symbols to each line. From the perspective of probability, the character-based language model is aiming for predicting the next character based on the context (previous two characters), in which start symbols and end symbols can be the forward reference and backward reference to make the trigram grammar a valid probability distribution.

```
1 def _preprocess_line(self, raw_line):
2     """
3     Convert raw lines to the required format". Requirements are as
4     ↪ follows:
5         1. Retain characters like English alphabets, spaces, digits, or '.'
6     ↪ characters exclusively.
7         2. Lowercase all characters
8         3. Convert all digits to 0
9     param:
10        raw_line (str): a sentence in its raw format
11    Output:
12        processed_line (str): the sentence in the required format
13    """
14    char_list = [char for char in raw_line]
15    processed_list = []
16    for char in char_list:
17        if re.match(r"[a-zA-Z]", char) or char.isspace():
18            processed_list.append(char.lower())
19        elif char.isdigit():
20            processed_list.append('0')
21        elif char == '.':
```

```

20         processed_list.append('.')
21     processed_line = ''.join(processed_list).replace('\n', '')
22
23     return processed_line
24
25 def _add_symbols(self, processed_line):
26     """
27     Add start and end symbols to processed lines. In other words, adding
    ↪ ## to the start and the end of previously processed line
28     param:
29         processed_line (str): a sentence in the required format but
    ↪ without symbols
30     Output:
31         complete_line (str): the sentence in the required format with
    ↪ symbols.
32     """
33     start_symbol, end_symbol = self.symbol, self.symbol
34     num_symbol = self.n_gram - 1
35     added_start_symbols, added_end_symbols = start_symbol * num_symbol,
    ↪ end_symbol * num_symbol
36
37     complete_line = added_start_symbols + processed_line +
    ↪ added_end_symbols
38
39     return complete_line
40

```

## 2 Examining a pre-trained model

In the given language model, we would assume that it has implemented Maximum Likelihood Estimation(MLE) along with add- $\alpha$  smoothing estimation. In the trigram distribution, we noticed loads of trigrams share the same estimated probability; for example, any character given the previous two characters are ‘#’, space, has the probability of 3.448e-02, that is, these trigrams share the same frequency in the corpus, which is quite rare in natural language data. Another counterintuitive observation is that  $P(space \mid \#\#)$  is not zero, whereas most space will not be the first character of a sentence.

Besides, we also conducted some experiments to narrow the range of  $\alpha$  based on the given training data. According to the perplexity shown in Section 5, there is a huge gap of perplexity between our implemented add- $\alpha$  (where  $\alpha = 0.8$ ) smoothing algorithm and the given model. We assumed it is a result of a large  $\alpha$ , since if the probability of each trigram is nearly the same, then the model will guess the third character according to a uniform distribution by giving the first two characters. Essentially, the model is not meaningful when the trigram distribution is uniform, as it does not learn anything but only takes stochastic guesses. To achieve the uniform distribution, the  $\alpha$  is supposed to be an extremely large number; therefore, the frequency of each trigram is approximately the same. Moreover, the value of  $\alpha$  should relate to the corpus size, whereas we do not have for given model thus fail to take it into consideration. By setting  $\alpha$  to  $10^1, 10^2, 10^3, 10^4, 10^5, 10^6, 10^7$ , the perplexities are 10.14, 16.17, 25.37, 29.28, 29.91, 29.98, 29.98. Following the increase of  $\alpha$ , the perplexity increases, which means the performance of the character-based language model becomes poorer. We narrowed the range of  $\alpha$  and found that when setting  $\alpha$  to 400, the perplexity of our implemented model (22.02) approximately equals that of given model (22.09). Hence, we guess the pre-trained model is trained through add- $\alpha$  smoothing estimation but with an inappropriately larger  $\alpha$  value.

### 3 Implementing a model: description and example probabilities (35 marks)

#### 3.1 Model description

We built a character-based trigram language model, estimating the probability distribution based on MLE and add- $\alpha$  smoothing. The combination of MLE and add- $\alpha$  smoothing is a straightforward and effective way deal with error handling by manual add  $\alpha$  to all possible trigrams in case that some trigrams, not shown in training data, cannot be identified by the language model. In addition, add- $\alpha$  smoothing can help on the problems of overfitting. According to MLE, our trigram parameter estimation will follow:

$$P(w_n|w_{n-1}, w_{n-2}) = \frac{C(w_{n-2}, w_{n-1}w_n)}{C(w_{n-2}w_{n-1})} \quad (1)$$

Therefore, estimated probability of a trigram equals the frequency that it appears in the corpus, divided by its history counts - the counts of trigrams sharing the same first two characters. MLE can well reflect the probability of high-frequency trigrams, whereas, if a trigram (in testing data) has never occurred in the training data, the method ignores the probability on the existence of the trigram, which is not the fact in most cases since training data could be limited or language keeps evolving. For solving this issue, we add a value  $\alpha$  (initially set to 1 and finally adjusted to 0.8, full details are in Appendix: Experimental Results) to each possible trigram (the size of possible trigrams is marked as  $V$  below) to smooth our results from MLE. After introducing a hyperparameter to our original estimation, each trigram parameter changed based on:

$$P(w_n|w_{n-1}, w_{n-2}) = \frac{C(w_{n-2}, w_{n-1}w_n) + \alpha}{C(w_{n-2}w_{n-1}) + V\alpha} \quad (2)$$

#### 3.2 Model excerpt

For potential characters after ‘ng’, we usually would expect nothing since ‘ng’ is usually part of the morpheme ‘ing’. If we consider ‘ng’ as one sound unit, which is usually phonemically realised as /ŋ/, it is always in a coda position and inhibited from appearing in an onset position in English. It is also not surprising that ‘ngf’, ‘ngt’, and ‘ngs’ are with higher probabilities because ‘ng’ does not have to be in the same morpheme or phoneme, which is rare(i.e., gingf, averdoingt). As for ‘ngs’, that could be a result of gerunds’ pluralisation, which is not common as well (i.e., belongings but not swimings).

```

1 {'ng ': 0.793155893536123, 'ng#': 0.0022813688212927796, 'ng.':
  ↳ 0.02636248415716101, 'ng0': 0.0010139416983523466, 'nga':
  ↳ 0.0035487959442332124, 'ngb': 0.0010139416983523466, 'ngc':
  ↳ 0.0010139416983523466, 'ngd': 0.004816223067173646, 'nge':
  ↳ 0.08593155893536136, 'ngf': 0.0022813688212927796, 'ngg':
  ↳ 0.0010139416983523466, 'ngh': 0.0010139416983523466, 'ngi':
  ↳ 0.0022813688212927796, 'ngj': 0.0010139416983523466, 'ngk':
  ↳ 0.0010139416983523466, 'ngl': 0.0035487959442332124, 'ngm':
  ↳ 0.0010139416983523466, 'ngn': 0.0022813688212927796, 'ngo':
  ↳ 0.007351077313054512, 'ngp': 0.0010139416983523466, 'ngq':
  ↳ 0.0010139416983523466, 'ngr': 0.012420785804816245, 'ngs':
  ↳ 0.021292775665399277, 'ngt': 0.013688212927756679, 'ngu':
  ↳ 0.0035487959442332124, 'ngv': 0.0010139416983523466, 'ngw':
  ↳ 0.0010139416983523466, 'ngx': 0.0010139416983523466, 'ngy':
  ↳ 0.0010139416983523466, 'ngz': 0.0010139416983523466}

```

## 4 Generating from models

The basic algorithm of this trigram model is predicting the next character conditioned on the previous two characters. Therefore, in our generating function, we applied a sliding window at a position to condition on the previous two characters, and then predict the character at the position by a non-deterministic approach based on a trigram probability distribution. Every time generating a character, the sliding window will move a step further to condition on the next two new characters. The non-deterministic approach is to randomly select the characters according to the probabilities. A trigram is more likely to be selected if it has a larger probability according to the distribution of the first two characters. However, this non-deterministic approach aims not to select the trigram with the highest probability but with a high rate, which determines the non-uniqueness of the generated text. Moreover, we need to set some proper rules for start symbols and end symbols for not invalidating the text. In our function, we do not count any start symbol or end symbol into the length. For start symbols, the generation of first character should always be conditioned on two starting symbols. It is not sensible to generate the end symbols when the text does not reach the length. Hence, this function will not consider generating end symbols by removing the probability of trigrams involving end symbols, and manually adding end symbols when halting.

```
1 'From trained English language model'
2 '##ing to ambpurese of to ound intions alf auuucras rate iver.fxjpfact of
   ↳ the res subts ou the betich evies.ysw.0zcpwcgeopolvent her sm parlies
   ↳ wo.dsdxmosidtes a poliarow hav ea plecg0bic of mrstion ampotimosten
   ↳ porecesis and spartureand the its it the offeed econ nominiell
   ↳ lobjecord thasto coadied ach##'
3
4 'From given English language model'
5 '##ye hose do you.aoes it.igh.u0yubye.hnwblos agong that ducks..zip itthat
   ↳ turt.00cgue.n you napplethin the mor.kjgc zippy on he thi.in it.sdphou
   ↳ le shund we dow thes are wany.gl you does thats the tere therefragook
   ↳ ats to you my bund.bbigh.ehpdvcpy.jmwrrouts at to you
   ↳ wat.wjhn.wmally.lice.t is.rrobbit##'
6
```

## 5 Computing perplexity

To evaluate how well our models on predicting the data, we should fit our models to dev sets and calculate the perplexity based on the equation below:

$$PP_M(\vec{w}) = P_M(w_1 \dots w_n)^{-\frac{1}{n}} \quad (3)$$

However, we do not have given dev sets to fit our models, then we have to fit them directly to the test data, which might have potential overfitting issues. Another issue is the simple product of each trigram's probability based on our model will be close to zero, which leads to some non-obvious results. Therefore, we will utilise cross-entropy, which is a logarithmic version of calculating perplexity and will give us an obvious number. All the perplexity is measured based on the add- $\alpha$  (where  $\alpha = 0.8$ ) smoothing.

$$PP_M(\vec{w}) = 2^{\log_2 P_M(w_1 \dots w_n)^{-\frac{1}{n}}} \quad (4)$$

Results:

```
1 'perplexity based on training.en': 8.858391703180724
2 'perplexity based on given model': 22.094457902881697
3 'perplexity based on training.de' : 23.408273543425153
4 'perplexity based on training.es' : 22.9957778613915
```

From the perplexity we gained, we can infer that our model based on 'training.en' is better at predicting the data set we fit in, because according to our equation, the higher the perplexity, the lower  $P_M(w_1 \dots w_n)$ . If we fit our English model to a new test data, and we get a higher perplexity, which means our model is not generalised enough to make predictions or the test data is noisy (perhaps not in English). However, this perplexity of our model is not sound enough, which probably results from the limited training data as only 1000 lines are involved in training data. Another reason could be the imperfect selection of smoothing algorithm as add- $\alpha$  is limited in many aspects.

## 6 Further improvements

In this section, we will seek further improvements based on some intuitions from the aspects of linguistics. Our core task is to conduct language identification. Therefore, firstly, we need to consider for a character-based model how languages are different in a ‘character’ sense, based on which we make further developments. Obviously, we have processed them to be the same in a semiotic sense; however, we could make provisional assumptions based on the random tri-characters, which provide helpful but subtle information for morphological parsing and word formation. Note that our proposes are some shallow aspects of morphological parsing, while the vision is much more complicated to handle. We will mainly take advantage of inflectional morphemes for two reasons: (i) Inflectional morphemes usually are less varied and have a strict selectional property (i.e., derivational affix like -ly can be attached to either a noun or a verb, whereas inflectional ones like -ing or -ed only go behind a verb.). (ii) Inflectional morphemes tend to have a broader usage, thus a higher frequency in any possible text.

Moreover, another idea might be valuable is whether there is a distinction clear enough to help us distinguish among different languages. Morpheme borrowing happens all the time, and the historical development which leads to language variation are also overlapping in the same language family (i.e., big vowel shift and umlaut in germanic languages → corresponding typological change). Despite some similarities, we can still point out the distinct morphemes if we did a morphological investigation on an individual language basis. Currently, we reckon that those annotations possibly will be done manually because of lack of understanding of cross-linguistic morphology and if we want to train a classifier to distinguish different morphemes that belong to different languages, that would be beyond the scale of this assignment. Then the basic idea of developments we could temporarily make is to highlight the language-specific morphemes, for example, we can get an average proportion of a language-specific morpheme, say, -ing from a English corpus(training data), and if that rate for another text is around or even higher, then it is likely to be in English. Similarly, we could use -ung in Geman identification task. However, there are exceptions to any rules, and these ‘distinct’ morphemes might happen to occur randomly (i.e., ung is also a part of English preterite stung) or in some unknown syntactic structure but at least what we proposed will increase the identification rate, and ideally in a larger corpus.

Word formation is another factor we think may give extra information about a language. However, most languages share the same word formation rules like compounding, blending and etc with the proportion differing a bit though. We believe this idea sheds light on morpheme extraction and benefit language identification on a morpheme-based model.



## Appendix: your code

Include a verbatim copy of your code for questions 1-5 here. If you answered question 6, you do *not* need to include that code.

Proprocessing:

```
1 import re
2 from math import log
3 from collections import defaultdict
4 import json
5 from random import choices
6
7
8 class Character_Model():
9
10     def __init__(self,
11                 training_data_path,
12                 distribution_write_path = '',
13                 write_distribution_status = False,
14                 smoothing_status = True,
15                 load_distribution_locally = False,
16                 distribution_file_path = 'trigram_distribution.json',
17                 symbol = '#',
18                 alpha = 0.8,
19                 n_gram = 3):
20         """
21         This class is for building a character based model for language
22         ↪ identification.
23         Params:
24             training_data_path (str): path of training data
25             distribution_write_path (str): path to write the probabilities
26             ↪ of trigrams
27             write_distribution_status (bool): whether to write
28             ↪ distribution
29             smoothing_status (bool): whether to apply add-alpha smoothing
30             load_distribution_locally (bool): for saving the time, after
31             ↪ generating the distribution file, it's not meaningful to compute the
32             ↪ distribution every time creating the object. Whether to load the
33             ↪ distribution from generated json file
34             distribution_file_path (str): path of a distribution
35             symbol (char): start and end symbol of lines
36             alpha (float): alpha value for add-alpha smooth
```

```

31         n_gram (int): here we set to 3 because this coursework is
↪ aiming to solve the problems of trigram
32
33         """
34         self.training_data_path = training_data_path
35         self.distribution_write_path = distribution_write_path
36         self.n_gram = n_gram
37         self.alpha = alpha
38         self.symbol = symbol
39
40         self.smoothing_status = smoothing_status
41
42         self.n_gram_frequency =
↪         self.get_n_gram_frequency(self.training_data_path)
43         if self.smoothing_status:
44             self.add_alpha_smoothing()
45
46         # load distribution from local file (computing distribution
↪ consumes a lot of time when smoothing algorithms are applied)
47         if load_distribution_locally:
48             with open(distribution_file_path, "r") as f:
49                 self.n_gram_distribution = json.load(f)
50             pass
51         else:
52             self.n_gram_distribution =
↪             self.get_n_gram_distribution(self.n_gram_frequency)
53
54         if write_distribution_status:
55             self.write_n_gram_distribution()
56
57
58     def _preprocess_line(self, raw_line):
59         """
60         Convert raw lines to the required format". Requirements are as
↪ follows:
61         1. Retain characters like English alphabets, spaces, digits, or '.'
↪ characters exclusively.
62         2. Lowercase all characters
63         3. Convert all digits to 0
64         param:
65         raw_line (str): a sentence in its raw format
66         Output:
67         processed_line (str): the sentence in the required format

```

```

68     """
69     char_list = [char for char in raw_line]
70     processed_list = []
71     for char in char_list:
72         if re.match(r"[a-zA-Z]", char) or char.isspace():
73             processed_list.append(char.lower())
74         elif char.isdigit():
75             processed_list.append('0')
76         elif char == '.':
77             processed_list.append('.')
78     processed_line = ''.join(processed_list).replace('\n', '')
79
80     return processed_line
81
82     def _add_symbols(self, processed_line):
83         """
84         Add start and end symbols to processed lines. In other words,
85         ↪ adding ## to the start and the end of processed line
86         param:
87             processed_line (str): sentence in the required format but
88             ↪ without symbols
89         Output:
90             complete_line (str): sentence in the required format with
91             ↪ symbols.
92
93         """
94         start_symbol, end_symbol = self.symbol, self.symbol
95         num_symbol = self.n_gram - 1
96         added_start_symbols, added_end_symbols = start_symbol *
97             ↪ num_symbol, end_symbol * num_symbol
98
99         complete_line = added_start_symbols + processed_line +
100             ↪ added_end_symbols
101
102         return complete_line
103
104     def get_n_gram_frequency(self, data_path):
105         """
106         1. Load the training data
107         2. Preprocess and add symbols on training data. Count on the
108         ↪ n-gram frequency
109         3. return the dict of n-gram frequencies

```

```

105         """
106
107         n_gram_frequency = defaultdict(int)
108
109         with open(data_path) as f:
110             for line in f:
111                 line = self._preprocess_line(line)
112                 line = self._add_symbols(line)
113
114                 for j in range(len(line) - self.n_gram):
115                     n_gram = line[j: j + self.n_gram]
116                     n_gram_frequency[n_gram] += 1
117
118         return n_gram_frequency
119
120     # specific trigram function
121     def add_alpha_smoothing(self):
122         """
123         This function can only applied to trigram model.
124         Apply add alpha smoothing to the language model:
125             1. get the set of characters
126             2. get all possible trgrams (removing some impossible cases
127 ↪ such as: 'a#a' '###')
128             3. Update the trigram frequency dict. Essentially, add alpha
129 ↪ smoothing is to change the frequency of each trigram.
130
131         Param:
132             self.symbol: start and end symbol of lines
133             self.alpha: the value of alpha added to frequency of every
134 ↪ trigram
135         """
136         # get the character set
137         char_vocab = []
138         for key in self.n_gram_frequency.keys():
139             char_vocab += set(list(key))
140             char_vocab = list(set(char_vocab))
141
142         trigram_vocab = []
143         # get all possibie trigrams
144         for i in range(len(char_vocab)):
145             for j in range(len(char_vocab)):
146                 for k in range(len(char_vocab)):
147                     trigram_vocab.append(f'{char_vocab[i]}'+

```

```

145         f'{char_vocab[j]}' +
146         f'{char_vocab[k]}')
147
148     # remove the impossible trigrams (patterns)
149     trigram_vocab.remove(self.symbol * 3)
150     for trigram in trigram_vocab:
151         if trigram[0] == self.symbol and trigram[2] == self.symbol:
152             trigram_vocab.remove(trigram)
153             continue
154
155         if trigram[1] == self.symbol and (trigram[0] != self.symbol
156     ↪ and trigram[2] != self.symbol):
157             trigram_vocab.remove(trigram)
158
159     # update the trigram frequency dict
160
161     for trigram in trigram_vocab:
162         if trigram in self.n_gram_frequency.keys():
163             self.n_gram_frequency[trigram] += self.alpha
164         else:
165             self.n_gram_frequency[trigram] = self.alpha
166
167     def _calculate_MLE(self, n_gram_name, frequency_dict):
168         """
169         Calculate the maximum likelihood estimation of a specific n_gram
170         ↪ based on the n_gram frequency
171         Params:
172             n_gram_name (str): name of n_gram
173             frequency_dict (dict): dictionary used to store the
174         ↪ frequencies of all n-grams
175         """
176         frequency = 0
177         n_gram_frequency = frequency_dict[n_gram_name]
178
179         # n_grams which shares the same n-1 characters
180         n_gram_pool = [n_gram for n_gram in frequency_dict.keys() if
181     ↪ n_gram.startswith(n_gram_name[0: self.n_gram - 1])]
182         for n_gram in n_gram_pool:
183             frequency += frequency_dict[n_gram]
184
185         mle = n_gram_frequency / frequency
186
187     return mle

```

```

184
185
186 def get_n_gram_distribution(self, n_gram_frequency):
187     """
188     Based on the n-gram frequency, calculating the probability of each
↪ n-gram
189     """
190     n_gram_distribution = {}
191     # need to consider the unseen trigram together with smoothing
192     for n_gram in sorted(n_gram_frequency.keys()):
193         mle = self._calculate_MLE(n_gram, n_gram_frequency)
194         n_gram_distribution[n_gram] = mle
195
196     return n_gram_distribution
197
198 def write_n_gram_distribution(self):
199     """
200     Write the n-gram distribution to a file in the JSON format
201     """
202     write_path = self.distribution_write_path
203     with open(f"{write_path}", "w") as outfile:
204
205         json.dump(self.n_gram_distribution, outfile)
206
207         print(f"file has been written at {write_path}")
208
209
210 # here trigram
211 def get_n_gram_probability(self, n_gram_head):
212     """
213     Get trigram probability based on the given first two characters
↪ (trigram head)
214     """
215     result = {}
216     n_gram_pool = [n_gram for n_gram in
↪ self.n_gram_distribution.keys() if
↪ n_gram.startswith(n_gram_head)]
217
218     for n_gram in n_gram_pool:
219         result[n_gram] = self.n_gram_distribution[n_gram]
220
221
222     return result

```

```

223
224
225 class Pre_trained_model():
226
227     def __init__(self,
228                 model_path = "./model-br.en"):
229         """
230         class for the given model
231         """
232         self.model_path = model_path
233         self.model_distribution = self.load_pre_trained_model()
234
235
236
237     def load_pre_trained_model(self):
238         """
239         Load the pre-trained model and get distribution of n-grams
240         """
241         n_grams, n_gram_distribution = [], {}
242
243         with open(self.model_path, 'r') as f:
244             for line in f:
245                 n_grams.append(re.sub(r"[\n\t]*", "", line))
246             for n_gram in n_grams:
247                 n_gram_distribution[n_gram[0: 3]] = float(n_gram[3: ])
248
249         return n_gram_distribution
250
251
252     def generate_from_LM(distribution, length = 300, start_symbols = '##',
253     ↪ end_symbols = '##'):
254         """
255         In this function, we manually set the start or end symbols don't count
256         ↪ on the length. In addition,
257         we manually add '##' to the start and then the generation of first
258         ↪ character depending on '##' and trigram distribution.
259         However, during generation, if the current length of a line doesn't
260         ↪ reach the limit, we remove all the possibilities of trigrams involving
261         ↪ '#'.
262
263         """
264         generated_txt = ''
265         # add the start symbols, if trigram => start symbols = '##'

```

```

261 generated_txt += start_symbols
262
263 while len(generated_txt) < length + 2:
264
265     trigram_head = generated_txt[-2: ]
266
267     trigram_pool, trigram_probability = [], []
268     for key, value in distribution.items():
269         if key.startswith(trigram_head):
270             if key[-1] != '#':
271                 trigram_pool.append(key);
272                 ↪ trigram_probability.append(value)
273
274     choice = choices(trigram_pool, trigram_probability)[0]
275     generated_txt += choice[-1]
276
277 generated_txt += end_symbols
278 return generated_txt
279
280 def read_test_file_by_lines(path):
281
282     with open(path) as f:
283         lines = f.readlines()
284         lines = [line.rstrip() for line in lines]
285
286     return lines
287
288 def preprocess_line(raw_line):
289
290     char_list = [char for char in raw_line]
291     processed_list = []
292     for char in char_list:
293         if re.match(r"[a-zA-Z]", char) or char.isspace():
294             processed_list.append(char.lower())
295         elif char.isdigit():
296             processed_list.append('0')
297         elif char == '.':
298             processed_list.append('.')
299     processed_line = ''.join(processed_list).replace('\n', '')
300
301     return processed_line
302

```



```

303 def preprocess_test_data(test_data_lines):
304
305     processed_data = []
306     for line in test_data_lines:
307         processed_line = '##' + preprocess_line(line) + '##'
308         processed_data.append(processed_line)
309
310     return processed_data
311
312 def calculate_perplexity(processed_data, distribution):
313     # could have some problems with the formula, need a further check
314     cross_entropy = 0
315     count = 0
316     for line in processed_data:
317         for i in range(len(line) - 3):
318             trigram = line[i: i + 3]
319             p = distribution[trigram]
320             minus_log_p = -1 * log(p, 2)
321             cross_entropy += minus_log_p
322             count += 1
323     cross_entropy_mean = cross_entropy / count
324     perplexity = pow(2, cross_entropy_mean)
325
326     return perplexity
327
328 def search_alpha(training_data_path, test_data_path, alpha_list):
329
330     raw_test_lines = read_test_file_by_lines(test_data_path)
331     processed_test_lines = preprocess_test_data(raw_test_lines)
332
333     result = {}
334     for alpha in alpha_list:
335         lm = Character_Model(training_data_path = training_data_path,
336                             load_distribution_locally=False,
337                             alpha=alpha,
338                             write_distribution_status=True,
339                             distribution_write_path=f"{alpha}.json")
340         perplexity = calculate_perplexity(processed_test_lines,
341                                         ↪ lm.n_gram_distribution)
342         result[str(alpha)] = perplexity
343
344     return result

```

```

345 # seems not valuable to build on command line based program; therefore, we
    ↳ choose to use the most straightforward way to run the file.
346
347 lm = Character_Model(training_data_path = "./code1st/training.es",
348                      distribution_write_path = './0.8.json',
349                      write_distribution_status = False,
350                      smoothing_status = True,
351                      load_distribution_locally = False,
352                      distribution_file_path='./0.8.json',
353                      symbol = '#',
354                      alpha = 0.8,
355                      n_gram = 3)
356
357 # model excerpt: distribution of ng
358 ng_distribution = lm.get_n_gram_probability('ng')
359 print(ng_distribution)
360
361 # txt generated from self-built language model
362 generated_txt_1 = generate_from_LM(length = 300, distribution =
    ↳ lm.n_gram_distribution)
363 print(generated_txt_1)
364 # txt generated from
365 pre_trained_lm = Pre_trained_model(model_path='./code1st/model-br.en')
366 generated_txt_2 = generate_from_LM(length = 300,
    ↳ distribution=pre_trained_lm.model_distribution)
367 print(generated_txt_2)
368
369 # calculate the perplexity of our model
370 test_lines = read_test_file_by_lines('./code1st/test')
371 processed_lines = preprocess_test_data(test_lines)
372 print(calculate_perplexity(processed_lines, lm.n_gram_distribution))
373
374 # # calculate the perplexity of given model
375 print(calculate_perplexity(processed_lines,
    ↳ pre_trained_lm.model_distribution))
376
377 # searching the alpha
378 result = search_alpha('./code1st/training.en', './code1st/test', [10, 100,
    ↳ 1000, 10000, 100000, 1000000])
379 print(result)

```

## Appendix: Experimental Results

This section shows the experimental results on different values of  $\alpha$  and the corresponding perplexities.

$\alpha$	Perplexity
0.2	8.97
0.4	8.87
0.6	8.86
0.8	8.86
1	8.87
1.2	8.89
1.4	8.91
1.6	8.93
1.8	8.96
2	8.99
10	10.14
$10^2$	16.17
400	22.02
$10^3$	25.37
$10^4$	29.29
$10^5$	29.91
$10^6$	29.98
$10^7$	29.98