

```
In [293]: 1 from numpy import array
2
3 import pandas as pd
4 import numpy as np
5 import re
6
7 import matplotlib.pyplot as plt
8 from sklearn.feature_extraction.text import CountVectorizer
9
10 START = "<s> "
11 STOP = "</s>"
12 UNK = "<UNK>"
```

```
In [318]: 1 # import libraries needed, read the dataset
2 import nltk, re, pprint, string
3 from nltk import word_tokenize, sent_tokenize
4
5 def add_sentence_tokens(sentences):
6     return ['{}{} {}'.format(START, s, STOP) for s in sentences]
7
8 def replace_tringletons(tokens):
9     vocab = nltk.FreqDist(tokens)
10    return [token if vocab[token] >= 3 else UNK for token in tokens]
11
12 def preprocess(sentences):
13     sentences = add_sentence_tokens(sentences)
14     tokens = ' '.join(sentences).split(' ')
15     tokens = replace_tringletons(tokens)
16     return tokens
17
```

```
In [359]: 1 # file = open('/Users/nehakardam/Documents/UWclasses /CSE NLP/A3/data_A
2 #file = open('/Users/nehakardam/Documents/UWclasses /CSE NLP/A3/data_A3
3 file = open('/Users/nehakardam/Documents/UWclasses /CSE NLP/A3/data_A3/
4 train = [l.strip() for l in file.readlines()]
5
```

```
In [361]: 1 unigram = preprocess(train,1)
2 bigrams = nltk.ngrams(unigram, 2)
3 trigrams = nltk.ngrams(unigram, 3)
```

```
In [362]: 1 freq_uni = nltk.FreqDist(unigram)
2 freq_bi = nltk.FreqDist(bigrams)
3 freq_tri = nltk.FreqDist(trigrams)
```

```
In [363]: len(freq_uni)
```

```
Out[363]: 9599
```

```
In [364]: 1 lm = LanguageModel(train, 1)
          2 print("Vocabulary size: {}".format(len(lm.vocab)))
          3 perplexity = lm.perplexity(train)
          4 print("Model perplexity: {:.3f}".format(perplexity))
```

Vocabulary size: 9599
Model perplexity: 537.947

```
In [367]: 1 lm = LanguageModel(train, 2)
          2 print("Vocabulary size: {}".format(len(lm.vocab)))
          3 perplexity = lm.perplexity(train)
          4 print("Model perplexity: {:.3f}".format(perplexity))
```

Vocabulary size: 9599
Model perplexity: 670.571

```
In [366]: 1 lm = LanguageModel(train, 3)
          2 print("Vocabulary size: {}".format(len(lm.vocab)))
          3 perplexity = lm.perplexity(train)
          4 print("Model perplexity: {:.3f}".format(perplexity))
```

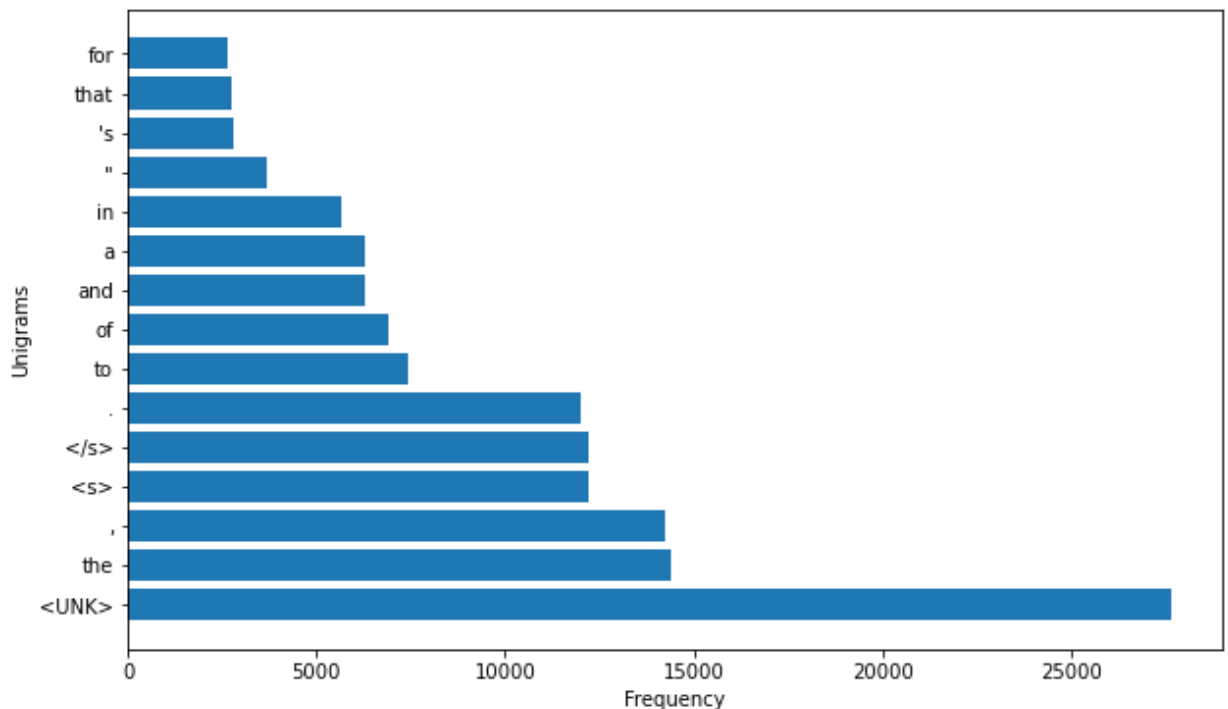
Vocabulary size: 9599
Model perplexity: 1592.478

```
In [357]: 1 # import math
          2
          3 # def compute_prob(freq, N):
          4 #     print ("processing")
          5 #     return { ngram: count / N for ngram, count in freq.items() }
          6
          7 # def perplexity(ngrams, freq):
          8 #     N = len(unigram)
          9 #     prob = compute_prob(freq, N)
         10 #     probabilities = [prob[ngram] for ngram in ngrams]
         11 #     print(probabilities)
         12 #     return math.exp((-1/N) * sum(map(math.log, probabilities)))
```

```

In [358]: 1 uni_k = list(freq_uni.keys())
          2 uni_v= list(freq_uni.values())
          3 bi_k = list(freq_bi.keys())
          4 bi_v= list(freq_bi.values())
          5 tri_k = list(freq_tri.keys())
          6 tri_v= list(freq_tri.values())
          7
          8 bi_k_str = ['_'.join(i) for i in bi_k]
          9 tri_k_str = ['_'.join(i) for i in tri_k]
         10 unigram_count_dict = dict(zip(uni_k,uni_v))
         11 bigram_count_dict = dict (zip(bi_k_str,bi_v))
         12 trigram_count_dict = dict (zip(tri_k_str,tri_v))
         13
         14 unigram_data_items = unigram_count_dict.items()
         15 unigram_data_list = list(unigram_data_items)
         16 unigram_df = pd.DataFrame(unigram_data_list,columns=[ 'Unigram', 'count' ])
         17 unigram_df = unigram_df.sort_values(by=[ 'count' ],ascending=False,ignore_
         18 plt.rcParams["figure.figsize"] = (10,6)
         19 y=unigram_df[ 'Unigram' ]
         20 x=unigram_df[ 'count' ]
         21 plt.barh(y, x)
         22 plt.ylabel( "Unigrams")
         23 plt.xlabel( "Frequency")
         24 # plt.title("Bigram Frequency Distribution")
         25 plt.show()

```



In [360]:

```

1  #!/bin/env python
2
3  import argparse
4  from itertools import product
5  import math
6  import nltk
7  from pathlib import Path
8
9  #!/bin/env python
10
11 import nltk
12
13 SOS = "<s> "
14 EOS = "</s>"
15 UNK = "<UNK>"
16
17 def add_sentence_tokens(sentences, n):
18     """Wrap each sentence in SOS and EOS tokens.
19
20     For n >= 2, n-1 SOS tokens are added, otherwise only one is added.
21
22     Args:
23         sentences (list of str): the sentences to wrap.
24         n (int): order of the n-gram model which will use these sentences.
25     Returns:
26         List of sentences with SOS and EOS tokens wrapped around them.
27
28     """
29     sos = SOS * (n-1) if n > 1 else SOS
30     return ['{}{} {}'.format(sos, s, EOS) for s in sentences]
31
32 def replace_singletons(tokens):
33     """Replace tokens which appear only once in the corpus with <UNK>.
34
35     Args:
36         tokens (list of str): the tokens comprising the corpus.
37     Returns:
38         The same list of tokens with each singleton replaced by <UNK>.
39
40     """
41     vocab = nltk.FreqDist(tokens)
42     return [token if vocab[token] >= 3 else UNK for token in tokens]
43
44 def preprocess(sentences, n):
45     """Add SOS/EOS/UNK tokens to given sentences and tokenize.
46
47     Args:
48         sentences (list of str): the sentences to preprocess.
49         n (int): order of the n-gram model which will use these sentences.
50     Returns:
51         The preprocessed sentences, tokenized by words.
52
53     """
54     sentences = add_sentence_tokens(sentences, n)
55     tokens = ' '.join(sentences).split(' ')
56     tokens = replace_singletons(tokens)

```

```

57     return tokens
58
59
60 def load_data(data_dir):
61     """Load train and test corpora from a directory.
62
63     Directory must contain two files: train.txt and test.txt.
64     Newlines will be stripped out.
65
66     Args:
67         data_dir (Path) -- pathlib.Path of the directory to use.
68
69     Returns:
70         The train and test sets, as lists of sentences.
71
72     """
73     train_path = data_dir.joinpath('train.txt').absolute().as_posix()
74     test_path = data_dir.joinpath('test.txt').absolute().as_posix()
75
76     with open(train_path, 'r') as f:
77         train = [l.strip() for l in f.readlines()]
78     with open(test_path, 'r') as f:
79         test = [l.strip() for l in f.readlines()]
80     return train, test
81
82
83 class LanguageModel(object):
84     """An n-gram language model trained on a given corpus.
85
86     For a given n and given training corpus, constructs an n-gram language
87     model for the corpus by:
88     1. preprocessing the corpus (adding SOS/EOS/UNK tokens)
89     2. calculating (smoothed) probabilities for each n-gram
90
91     Also contains methods for calculating the perplexity of the model
92     against another corpus, and for generating sentences.
93
94     Args:
95         train_data (list of str): list of sentences comprising the training
96         n (int): the order of language model to build (i.e. 1 for unigram)
97         laplace (int): lambda multiplier to use for laplace smoothing
98
99     """
100
101     def __init__(self, train_data, n, laplace=1):
102         self.n = n
103         self.laplace = laplace
104         self.tokens = preprocess(train_data, n)
105         self.vocab = nltk.FreqDist(self.tokens)
106         self.model = self._create_model()
107         self.masks = list(reversed(list(product((0,1), repeat=n))))
108
109     def _smooth(self):
110         """Apply Laplace smoothing to n-gram frequency distribution.
111
112         Here, n_grams refers to the n-grams of the tokens in the training
113         corpus, while m_grams refers to the first (n-1) tokens of each n-gram.

```

```

114
115     Returns:
116         dict: Mapping of each n-gram (tuple of str) to its Laplace-
117             probability (float).
118
119     """
120     vocab_size = len(self.vocab)
121
122     n_grams = nltk.ngrams(self.tokens, self.n)
123     n_vocab = nltk.FreqDist(n_grams)
124
125     m_grams = nltk.ngrams(self.tokens, self.n-1)
126     m_vocab = nltk.FreqDist(m_grams)
127
128     def smoothed_count(n_gram, n_count):
129         m_gram = n_gram[:-1]
130         m_count = m_vocab[m_gram]
131         return (n_count + self.laplace) / (m_count + self.laplace)
132
133     return { n_gram: smoothed_count(n_gram, count) for n_gram, count in n_grams.items() }
134
135 def _create_model(self):
136     """Create a probability distribution for the vocabulary of the
137
138     If building a unigram model, the probabilities are simple relative
139     of each token with the entire corpus.
140
141     Otherwise, the probabilities are Laplace-smoothed relative frequencies.
142
143     Returns:
144         A dict mapping each n-gram (tuple of str) to its probability.
145
146     """
147     if self.n == 1:
148         num_tokens = len(self.tokens)
149         return { (unigram,): count / num_tokens for unigram, count in self.n_grams.items() }
150     else:
151         return self._smooth()
152
153 def _convert_oov(self, ngram):
154     """Convert, if necessary, a given n-gram to one which is known
155
156     Starting with the unmodified ngram, check each possible permutation
157     with each index of the n-gram containing either the original token
158     when the model contains an entry for that permutation.
159
160     This is achieved by creating a 'bitmask' for the n-gram tuple,
161     each flagged token for <UNK>. Thus, in the worst case, this function
162     checks all possible n-grams before returning.
163
164     Returns:
165         The n-gram with <UNK> tokens in certain positions such that
166         contains an entry for it.
167
168     """
169     mask = lambda ngram, bitmask: tuple((token if flag == 1 else "<UNK>") for flag, token in zip(bitmask, ngram))
170

```

```

171         ngram = (ngram,) if type(ngram) is str else ngram
172         for possible_known in [mask(ngram, bitmask) for bitmask in self.model.itertools.combinations(
173             self.model.get_vocab().keys(), ngram[0]):
174             if possible_known in self.model:
175                 return possible_known
176
177     def perplexity(self, test_data):
178         """Calculate the perplexity of the model against a given test data.
179
180         Args:
181             test_data (list of str): sentences comprising the training data.
182         Returns:
183             The perplexity of the model as a float.
184
185         """
186         test_tokens = preprocess(test_data, self.n)
187         test_ngrams = nltk.ngrams(test_tokens, self.n)
188         N = len(test_tokens)
189
190         known_ngrams = (self._convert_oov(ngram) for ngram in test_ngrams)
191         probabilities = [self.model[ngram] for ngram in known_ngrams]
192
193         return math.exp((-1/N) * sum(map(math.log, probabilities)))
194
195     def _best_candidate(self, prev, i, without=[]):
196         """Choose the most likely next token given the previous (n-1) tokens.
197
198         If selecting the first word of the sentence (after the SOS token),
199         the i'th best candidate will be selected, to create variety.
200         If no candidates are found, the EOS token is returned with probability 1.
201
202         Args:
203             prev (tuple of str): the previous n-1 tokens of the sentence.
204             i (int): which candidate to select if not the most probable.
205             without (list of str): tokens to exclude from the candidate list.
206         Returns:
207             A tuple with the next most probable token and its corresponding probability.
208
209         """
210         blacklist = ["<UNK>"] + without
211         candidates = ((ngram[-1], prob) for ngram, prob in self.model.itertools.combinations(
212             self.model.get_vocab().keys(), n))
213         candidates = filter(lambda candidate: candidate[0] not in blacklist, candidates)
214         candidates = sorted(candidates, key=lambda candidate: candidate[1])
215         if len(candidates) == 0:
216             return ("</s>", 1)
217         else:
218             return candidates[0 if prev != () and prev[-1] != "<s>" else i]
219
220     def generate_sentences(self, num, min_len=12, max_len=24):
221         """Generate num random sentences using the language model.
222
223         Sentences always begin with the SOS token and end with the EOS token.
224         While unigram model sentences will only exclude the UNK token,
225         this model will also exclude all other words already in the sentence.
226
227         Args:
228             num (int): the number of sentences to generate.
229             min_len (int): minimum allowed sentence length.

```



```

228         max_len (int): maximum allowed sentence length.
229     Yields:
230         A tuple with the generated sentence and the combined probal
231         (in log-space) of all of its n-grams.
232
233     """
234     for i in range(num):
235         sent, prob = ["<s>"] * max(1, self.n-1), 1
236         while sent[-1] != "</s>":
237             prev = () if self.n == 1 else tuple(sent[-(self.n-1):])
238             blacklist = sent + (["</s>"] if len(sent) < min_len else [])
239             next_token, next_prob = self._best_candidate(prev, i, v, blacklist)
240             sent.append(next_token)
241             prob *= next_prob
242
243             if len(sent) >= max_len:
244                 sent.append("</s>")
245
246         yield ' '.join(sent), -1/math.log(prob)
247

```

Linear Interpolation Smoothing

Reference: <https://github.com/ErolOZKAN-/Language-Modelling/blob/master/src/ngram.py>.
<https://github.com/ErolOZKAN-/Language-Modelling/blob/master/src/ngram.py>.

```

In [379]: 1 def ppl_interpolation(sentences, lambda_set):
2         '''Returns the interpolated perplexity of the given list of sentences'''
3         sentences = unigram
4         counter = 0
5         tmp = 0
6         for i in range(len(sentences)):
7             sentence = sentences[i]
8             for j in range(len(sentence) - 2):
9                 trigram_unit = (sentence[j], sentence[j + 1], sentence[j + 2])
10                if trigram_unit in trigram_prob:
11                    tri_gram_prob = trigram_prob[trigram_unit]
12                else:
13                    tri_gram_prob = 0
14
15                bigram_unit = (sentence[j], sentence[j + 1])
16                if bigram_unit in bigram_prob:
17                    bigram_prob = bigram_prob[bigram_unit]
18                else:
19                    bigram_prob = 0
20
21                # unigram_prob = one_gram_add_one_prob.get(sentence[j], 0)
22
23                prob = lambda_set[0] * unigram_prob + lambda_set[1] * bigram_prob + lambda_set[2] * tri_gram_prob
24                if prob == 0:
25                    tmp += 0
26                else:
27                    tmp += math.log(prob, 2)
28                counter += 1
29            # entropy = prob of each token / number of tokens
30            entropy = -1 / counter * tmp
31            perplexity = math.pow(2, entropy)
32        return perplexity

```

```
In [380]: 1 print("\n-----INTERPOLATION-----")
2 lambda_set = []
3 lambda_set.append([0.5, 0.3, 0.2])
4 lambda_set.append([0.8, 0.1, 0.1])
5 lambda_set.append([0.1, 0.8, 0.1])
6 lambda_set.append([0.1, 0.1, 0.8])
7 lambda_set.append([0.6, 0.2, 0.2])
8 lambda_set.append([0.2, 0.6, 0.2])
9 lambda_set.append([0.2, 0.2, 0.6])
10 lambda_set.append([0.4, 0.3, 0.3])
11 lambda_set.append([0.3, 0.4, 0.3])
12 lambda_set.append([0.3, 0.3, 0.4])
13 lambda_set.append([0.2, 0.4, 0.4])
14 lambda_set.append([0.4, 0.2, 0.4])
15 lambda_set.append([0.4, 0.4, 0.2])
16 lambda_set.append([0.1, 0.4, 0.5])
17 lambda_set.append([0.1, 0.3, 0.6])
18 lambda_set.append([0.1, 0.2, 0.7])
19 lambda_set.append([0.05, 0.15, 0.8])
20 lambda_set.append([0.05, 0.05, 0.9])
21 for s in lambda_set:
22     print("Lambda_Set: ", lambda_set)
23     ppl_score = ppl_interpolation(train, s)
24     print("perplexity score of test data:", ppl_score)
```

-----INTERPOLATION-----

```
Lambda_Set: [[0.5, 0.3, 0.2], [0.8, 0.1, 0.1], [0.1, 0.8, 0.1], [0.1, 0.1, 0.8], [0.6, 0.2, 0.2], [0.2, 0.6, 0.2], [0.2, 0.2, 0.6], [0.4, 0.3, 0.3], [0.3, 0.4, 0.3], [0.3, 0.3, 0.4], [0.2, 0.4, 0.4], [0.4, 0.2, 0.4], [0.4, 0.4, 0.2], [0.1, 0.4, 0.5], [0.1, 0.3, 0.6], [0.1, 0.2, 0.7], [0.05, 0.15, 0.8], [0.05, 0.05, 0.9]]
```

```
-----
-
NameError                                Traceback (most recent call last)
<ipython-input-380-63efe8819287> in <module>
    21 for s in lambda_set:
    22     print("Lambda_Set: ", lambda_set)
--> 23     ppl_score = ppl_interpolation(train, s)
    24     print("perplexity score of test data:", ppl_score)

<ipython-input-379-fa79245e1b47> in ppl_interpolation(sentences, lambda_set)
     8         for j in range(len(sentence) - 2):
     9             trigram_unit = (sentence[j], sentence[j + 1], sentence[j + 2])
--> 10             if trigram_unit in trigram_prob:
    11                 tri_gram_prob = trigram_prob[trigram_unit]
    12             else:

NameError: name 'trigram_prob' is not defined
```

In []:

1