# Accelerated Natural Language Processing
# Assignment 1
# Report

s1882930

s1842611

# Q1.

Here are the codes:

```
1.    # judge the character belongs to the alphabet
2.    def is_alphabet(char):
3.        if (char >= 'a' and char <= 'z'
4.            or
5.            char >= 'A' and char <= 'Z'
6.        ):
7.            return True
8.        else:
9.            return False
10.
11.   def preprocess_line(a):
12.       b = ['##']
13.   #set  '##'as start of a sentence
14.   #in order to get the probability of first character when we use bigram
15.   #and trigram model
16.       for i in range(len(a)):
17.           if a[i].isdigit():
18.   #get the digit of the sentence and change it into '0'
19.               b.append('0')
20.           if a[i] == " " or a[i] == '.' or is_alphabet(a[i]):
21.   #judge it if is the target characters
22.               b.append(a[i])
23.       b.append('#')
24.   #set  '#' as end of a sentence
25.   #in order to calculate the probability of the trigrams
26.   #it will be convenient when computing perplexity
27.       c = ''.join(b).lower().
28.   return c
```

# Q2.

I guess that it uses maximum likelihood estimation and add-one smoothing, when calculating the probabilities of trigrams. Reasons are following.

Firstly, from the distribution of the data given by **model_br.en**, there are 26100 different trigrams in the model. The probabilities range from 1.021e-5 to 0.9986, which means that there is no zero probability. Secondly, given the previous same two history words, the sum probabilities of different trigram is one. Thirdly, the mode(the most frequent probability value) value of all data is 0.0333 and it occurs 12632 times. Besides, 0.0333 is very important because it approximates 1/30, whose denominator is just the type number (the vocabulary size) of each trigram model. According to the mode, I assume that the language model uses the add-one smoothing method to get the probability

of each trigram.

## Q3.

We use interpolation to build up the language model with **training.en** data, mixing the probability estimates from all the n-gram estimators, weighing and combining the trigram, bigram, and unigram. Because when there is only a small training set, using the ordinary equation

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

to get the trigram probability will generate many zero probabilities. Besides, the counts of two previous characters (bigram) may also be zero, which means the probability value is infinite.

Compared with add-one smoothing and add-$\alpha$ smoothing which distribute the same probability to each unknown trigram which doesn't appear in the given corpus, the interpolation method can distribute different probabilities to each unknown trigram according to the bigram model and unigram model.

In order to solve it, we use interpolation method to set up language model.

$$P(w_n|w_{n-2}w_{n-1}) = \lambda_1 P(w_n|w_{n-2}w_{n-1}) + \lambda_2 P(w_n|w_{n-1}) + \lambda_3 P(w_n)$$
$$where \; \lambda_1 + \lambda_2 + \lambda_3 = 1$$
$$0 < \lambda_i < 1$$

Actually, each $\lambda$ weight is computed by conditioning on the context. But considering that we have one small training set, we make the assumption on that each $\lambda$ makes the same contribution to the language model, and set $\lambda_1 = \lambda_2 = \lambda_3 = 1/3$.

Codes are following:

```
1.   # get the ngram model and compute the counts of n_grams from the training set
2.   #line_counts is the total number of the whole passage
3.   def ngram_counts(n, file):
4.       counts = defaultdict(int)
5.       line_count = 0
6.       with open(file) as f:
7.           for line in f:
8.               line = preprocess_line(line)
9.               line_count += 1
10.              if len(line) > n:
11.                  for i in range(3 - n, len(line) - (n - 1), 1):
12.                      gram = line[i:i + n]
13.  # ngram model counts
14.                      counts[gram] += 1
15.  return counts, line_count
16.
17.  # calculate the unigram probability
18.  #the input is the counts of unigrams in the training set
19.  def u_p(origin_counts):
20.      p_model = defaultdict(float)
```

```python
21.     tokens = 0
22.     types = 0
23.     for i in origin_counts:
24.         tokens += origin_counts[i]
25.     for i in origin_counts:
26.         p_model[i] = origin_counts[i] / tokens
27.     types = len(origin_counts)
28.     return p_model
29.
30.
31. #uc is an dictionary contains the the counts of unigrams and probabilities
32. #origin_b_c is the counts of bigrams in the training set
33. #return the the whole dictionary containing all bigrams
34. def b_model_cor(uc, origin_b_c):
35.     bmc = defaultdict(int)
36.     for i in uc:
37.         for j in uc:
38.             temp = []
39.             temp.append(i)
40.             temp.append(j)
41.             word = ''.join(map(str, temp))
42.             if (word in origin_b_c.keys()):
43.                 bmc[word] = origin_b_c[word]
44.             else:
45. #if the previous characters donnot appear in the passage, it is assigned to 0
46.                 bmc[word] = 0
47.     del bmc['##']
48.     return bmc
49.
50. # get  the counts of trigrams
51. # tc is the original counts of the trigram in the training set
52. def t_model_cor(uc, tc):
53.     tmc = defaultdict(int)
54.     for i in uc:
55.         for j in uc:
56.             for m in uc:
57.                 temp = []
58.                 temp.append(i)
59.                 temp.append(j)
60.                 temp.append(m)
61.                 word = ''.join(map(str, temp))
62.                 if (word in tc.keys()):
63.                     tmc[word] = tc[word]
64.                 else:
```

```python
65.                 tmc[word] = 0
66.     #delete the impossible trigrams
67.     for i in uc:
68.         for j in uc:
69.             for m in uc:
70.                 if (i == '#' and m == '#'):
71.                     del tmc[i + j + m]
72.     for i in uc:
73.         for j in uc:
74.             for m in uc:
75.                 if (i != '#'):
76.                     if (j == '#'):
77.                         del tmc[i + j + m]
78.     return tmc
79.
80.     #   bcp is a whole counts for bigram dictionary
81.     def b_gram(bcp, uc):
82.         p_model = defaultdict(float)
83.         tokens = 0
84.         types = len(bcp)
85.         for i in bcp:
86.             tokens += bcp[i]
87.         for i in bcp:
88.             p_model[i] = bcp[i] / uc[i[0]]
89.         return p_model
90.
91.     #set up the trigram model contain the zero probabilities
92.     #tcp is a dictionary containing the probability of the trigrams
93.     #bcp is a dictionary containing the probability of the bigrams
94.     def t_gram(tcp, bcp):
95.         p_model = defaultdict(float)
96.         tokens = 0
97.         types = len(tcp)
98.         line_number = 0
99.         for i in tcp:
100.            tokens += tcp[i]
101.        for i in tcp:
102.            if (i[:2] == '##'):
103.                line_number = line_number + tcor_c[i]  # calculate the line number
104.        for i in tcp:
105.            if (bcp[i[:2]] != 0):
106.                p_model[i] = tcp[i] / bcp[i[:2]]
107.            else:
108.                if (i[:2] == '##'):
```

```python
109. # '##'represents the start of a sentence
110.             p_model[i] = tcp[i] / line_number
111.         else:
112.             p_model[i] = 0
113.     return p_model
114.
115. # using the a1,a2,a3 to set up the interpolation
116. def interplotaion(u, b, t, a1, a2, a3):
117.     p_model = defaultdict(float)
118.     for i in t:
119.         p_model[i] = a1 * t[i] + a2 * b[i[1:]] + a3 * u[i[2:]]
120.     return p_model
121. #-------main body--------------------
122.
123. # generate the basic counts dictionary
124. u_counts, line_count = ngram_counts(1, "/Users/chenwenjun/PycharmProjects/w2/venv/training.en")
125. b_counts, line_count = ngram_counts(2, "/Users/chenwenjun/PycharmProjects/w2/venv/training.en")
126. t_counts, line_count = ngram_counts(3, "/Users/chenwenjun/PycharmProjects/w2/venv/training.en")
127. #generate the unigram model
128. ucor_c = defaultdict(int)
129. ucor_c = u_counts
130.
131. # generate unigram model
132. ucor_p = defaultdict(float)
133. ucor_p = u_p(ucor_c)
134.
135.
136. # generate bigram model
137. bcor_c = defaultdict(int)
138. bcor_c = b_model_cor(u_counts, b_counts)
139. bcor_p = defaultdict(float)
140. bcor_p = b_gram(bcor_c, u_counts)
141.
142. # generate trigram model
143. tcor_c = defaultdict(int)
144. tcor_c = t_model_cor(u_counts, t_counts)
145. tcor_p = defaultdict(float)
146. tcor_p = t_gram(tcor_c, bcor_c)
147.
148. # generate interpolation
149. f1 = defaultdict(float)
150. f1 = interplotaion(ucor_p, bcor_p, tcor_p,1/3, 1/3, 1/3)
```

Here is the dictionary of **ng** :

| | |
|---|---|
| 'ng#': | 0.0027488455108938257 |
| 'ng.': | 0.014032523221172578 |
| 'ng ': | 0.428113494670319 |
| 'ng0': | 0.0015301686809280684 |
| 'nga': | 0.03763821654873327 |
| 'ngb': | 0.0034305047502100434 |
| 'ngc': | 0.008716150714147225 |
| 'ngd': | 0.011948246928464492 |
| 'nge': | 0.10665589096829017 |
| 'ngf': | 0.006691569056128392 |
| 'ngg': | 0.0065382569120528596, |
| 'ngh': | 0.032971380415634334 |
| 'ngi': | 0.0696990250537054 |
| 'ngj': | 0.00042612292380275323 |
| 'ngk': | 0.0013515413946875202, |
| 'ngl': | 0.012165034364289826 |
| 'ngm': | 0.008158747495637563 |
| 'ngn': | 0.02429830303543762 |
| 'ngo': | 0.042925607838841084 |
| 'ngp': | 0.0072139598009435795 |
| 'ngq': | 0.0003271246687778711, |
| 'ngr': | 0.06909447177113728 |
| 'ngs': | 0.027610229783556043, |
| 'ngt': | 0.032780557480811844 |
| 'ngu': | 0.027213327243800767 |
| 'ngv': | 0.0028817100756155889, |
| 'ngw': | 0.004304271957603568, |
| 'ngx': | 0.0005552510825308602 |
| 'ngy': | 0.006884771668779344 |
| 'ngz': | 0.00009469398306727848 |

We expect that the sum of all conditional probability value should be one, and the real sum of all probabilities is 0.9986. Because when we pin down the first two characters, the sum of all conditional probabilities should be one. The reason why there is a little differences is that computer uses float variable to store the probability values, and it will lost accuracy.

Overall, the probability values extracted from the language model we build match what we expect.

## Q4.

Firstly we create a possible generated letter and save them as a list, which removes the symbol '#', because once generated # means that is the end of the characters generation, so removing the '#' is to ensure the number of characters generated, and then set the initial character '# #'. Next, we set

two empty lists to store the different trigram types and corresponding probability, and then extract the last two characters of the currently generated string, searching for the corresponding different type of trigram from our established language model and storing corresponding probabilities in the previously set empty list, and then store the probability value in the form of an array to facilitate subsequent selection by probability. Then, since we remove the case of generating '#', the sum value of the different type trigrams with given the same two history characters is not 1, and the function of the latter **random.choice** cannot be used, so normalization is performed, and each trigram probability value is divided by the sum of the probability values of all such trigrams. This doesn't affect the generation of the next character, because the probability values of all cases are expanded by the same times, and then we use the **random.choice** function with the previously extracted probability value to generate the corresponding character. Finally, the generated character is connected with the previously generated string, and the loop is repeated until the desired number of characters are generated.

```
1.   def generate_from_LM(distribution,N):
2.       y=['r','e','s','u','m','p','t','i','o','n',' ',
3.          'f','h','d','c','l','a','j','y','0','b','w','k','g','v','.','q','x','z']
4.       generate=[]
5.       w='##'
6.       for j in range(N):
7.           outcomes=[]
8.           probs=[]
9.           for i in y :
10.              a=[w[-2:]]
11.              a.append(i)
12.              a=''.join(a)
13.              outcomes.append(a)
14.              probs.append(distribution.get(a))
15.          probs = np.array(probs)
16.          summ=np.sum(probs)
17.          probs=probs/summ
18.          index = np.random.choice(outcomes, p = probs.ravel())
19.          generate.append(index)
20.          w=w+generate[-1][-1:]
21.      return print(w)
22.  generate_from_LM(f1,300)
23.  generate_from_LM(tr_frombr,300)
```

sequence generated from the model trained from "training.en":

   **##fartiew ht ha of epa tontosindonov. tihow prin emowor    nne reievere gucysshefiisshi ii inelide ds meral mpetety f f at    c daiven tly amif blar exemrtteptisin t do teunthan a encecl this sut inrt.albe wme of mldidider tme icoralt un orolighelheso buracot t iilite po weangomenfipsew ein ahatte t tunt**

sequence generated from the model trained from "mode-br.en":

**##bou dont oner.bcwn.auld ou theres.his it there can his is fore some.oh ah bood i samel a like falk.iyuirdalike thers. done.turn thats andy an to you like yought    at.ozgvfnight.h0zy dfkhqhcphat.ekaboy.srbwair dons.vqs the this gook to pre thwuoats to thi ch hem.mxize dtie.gcqn.y hing way now way opere**

We noticed that the sequence generated more '.' character based on teacher's model than our own model. And there are more real English words generated in teacher's model. Because the teacher's model is established on a much larger corpus than mine which is just based on an article. The teacher's model has more sentences, so simultaneously the model has more possibility of the character followed by '.'. Besides much larger corpus means more precise probability, so based on the teacher's model we can generate word better.

# Q5.

Followings are about the function about how to calculate the perplexity with the model generated from the "training.en"document

```
1.   def perplexicity(file, n, model):
2.       counts = defaultdict(int)
3.       #set up a dictionary which contains the probabilities and the trigrams
4.       with open(file) as f:
5.           for line in f:
6.               line = preprocess_line(line)
7.               if len(line) > n:
8.                   for i in range(3 - n, len(line) - (n - 1), 1):
9.                       gram = line[i:i + n]
10.                      if (model[gram] > 0):
11.                          counts[gram] = -log(model[gram], 2)
12.  #compute the log and add them together
13.                      else:
14.                          counts[gram] = 0
15.     # ngram model counts
16.     return  2**(sum(counts.values()) / len(counts)  )
```

When we input our language model and the test set document, it is fast to get the perplexity value.

```
1.   print(perplexity('/Users/chenwenjun/PycharmProjects/w2/venv/test',3,f1))
2.   print(perplexity('/Users/chenwenjun/PycharmProjects/w2/venv/training.es',3,f1))
3.   print(perplexity('/Users/chenwenjun/PycharmProjects/w2/venv/training.de',3,f1))
```

Here are the outputs :
17.933215118487485
26.46348580088351

31.686271894350266

From the output, we can see that the perplexity of English article is the lowest just the same as what we expect, because we train it from the English training set. But if we use other file to calculate the perplexity and get a lower result, it is unreasonable to prove that it belongs to English. Because all the sentences have to be processed by '**preprocess_line**' function, and in this circumstance, many characters have been lost. For example, if we use a paragraph which contains Chinese, English, and Japanese, the '**preprocess_line**' function will filter the Chinese and Japanese. So we can't only use perplexity as evidence to judge the type of language.

# Q6.

In the previous program, we assume that unigram, bigram, and trigram have the same contribution to the language model. Under this assumption, we set the parameters all the same. However, if we want to find the better parameters( $\lambda_1, \lambda_2, \lambda_3$ ) to build up a stable and robust models, it is evitable to use other methods to solve this problem.

Based on this question, we try to find the best parameters to set up a better language model with lower perplexity.

At first, we use loop to find better parameters:

```
1.   #min is the lowest perplexity
2.   #min_i,min_j,min_m represent the parameter λ1,λ2,λ3
3.   min=10000
4.   min_i=min_j=min_m=0
5.   for i in np.arange(0.1,0.9,0.01):
6.       for j in np.arange(0.1, 0.9, 0.01):
7.           if(i+j<1):
8.               m=1-i-j
9.               f1 = interplotaion(ucor_p, bcor_p, tcor_p, i, j, m)
10.              per=pe('/Users/chenwenjun/PycharmProjects/w2/venv/validation_set', 3, f1)
11.              if(per<min):
12.                  min=per
13.                  min_i=i
14.                  min_j=j
15.                  min_m = m
```

This function wants to find the hyperparameters within a defined step size and returns the minimum perplexity value. Besides, the right selection of smaller step sizes boosts the accuracy of the perplexity value returned by the algorithm.

We extract a validation set using 20% the original training set 'training.en'. At the same time, we set the step size 0.01 and let the computer do the loop until find the best parameters. Finally, we can get the parameters: 0.72 for trigram( $\lambda_1$ ), 0.15 for bigram( $\lambda_2$ ), and 0.13 for unigram( $\lambda_3$ ).

Then we use the new parameters to test the other languages and find the differences are more obvious.

```
1.  f2=interplotaion(ucor_p, bcor_p, tcor_p, 0.72, 0.15, 0.13)
2.  print(perplexity('/Users/chenwenjun/PycharmProjects/w2/venv/test', 3, f2))
3.  print(perplexity('/Users/chenwenjun/PycharmProjects/w2/venv/training.es', 3,
       f2))
4.  print(perplexity('/Users/chenwenjun/PycharmProjects/w2/venv/training.de', 3,
       f2))
```

Here are the results:

18.888578304183984

34.784594585819754

46.416519514021324

We think if we can have some extra data that is different from training data and test data. We can turn this data into development data. We define c'(u, v, w) to represent the parameters that the c' (u, v, w) appears in the development dataset. It is easy to calculate the log-likelihood values on the development dataset. The specific formula is as follows:

$$
\begin{aligned}
L(\lambda_1, \lambda_2, \lambda_3) &= \sum_{u,v,w} c'(u, v, w) \log q(w|u, v) \\
&= \sum_{u,v,w} c'(u, v, w) \log \left( \lambda_1 \times q_{ML}(w|u, v) + \lambda_2 \times q_{ML}(w|v) + \lambda_3 \times q_{ML}(w) \right)
\end{aligned}
$$

We choose the value of $\lambda$ so that the larger the value of $L(\lambda_1, \lambda_2, \lambda_3)$ is, the better outcome we can get. Therefore, the value of $\lambda$ is as follows:

$$
\arg \max_{\lambda_1, \lambda_2, \lambda_3} L(\lambda_1, \lambda_2, \lambda_3)
$$