# CS 458 - Natural Language Processing
# Assignment 01
# Batool Ali Akbar - ba07612

1. (40 points) Resolving Ambiguities Between DD/MM/YYYY and MM/DD/YYYY Date Formats

   (a) (10 points) Python Script

   **Solution:**

```python
import re

us_spelling = set([
    "color", "honor", "favorite", "theater", "center", "meter",
    "defense", "license", "analyze", "realize", "organize",
    "traveler", "jewelry", "program", "aluminum", "plow", "catalog",
    "gray", "leukemia", "tire",  "fall", "thanksgiving", "memorial day",
    "labor day", "fourth of july", "veterans day", "super bowl",
    "nba finals", "world series", "washington dc", "new york", "california",
    "chicago", "state", "governor", "semester","elementary school",
    "high school", "college", "fiscal year"
])

british_spelling = set([
    "colour", "honour", "favourite", "theatre", "centre", "metre",
    "defence", "licence", "analyse", "realise", "organise",
    "traveller", "jewellery", "programme", "aluminium", "plough",
    "catalogue", "grey", "leukaemia", "tyre",  "autumn", "christmas",
    "boxing day", "bank holiday", "easter", "good friday", "new year's day",
    "bonfire night", "summer holidays", "half term", "london", "paris",
    "berlin", "european union", "british", "england", "united kingdom",
    "favourite", "theatre", "university", "headteacher", "primary school",
    "secondary school", "european parliament", "fifa world cup", "uefa",
    "euros", "football"
])

corpus = open("date_format_dd_mm_yyyy.txt", "r")
corpus = corpus.read()
corpus = corpus.lower()
corpus_set = corpus.split()
corpus_set = set(corpus_set)

output_file = open("q1_batool_ba07612.txt", "w", encoding='utf-8')

dates = re.findall(r'\d{1,2}\/\d{1,2}\/\d{1,4}', corpus)

for d in dates:
    if int(d[:2]) > 12:
        format = "DD/MM/YYYY"
    elif int(d[3:5]) > 12:
        format  = "MM/DD/YYYY"
    else:
        if len(corpus_set.intersection(us_spelling)) != 0:
            format = "MM/DD/YYYY"
        elif len(corpus_set.intersection(british_spelling)) != 0:
            format = "DD/MM/YYYY"
        else:
            format = "ambiguous"
    output_file.write(d + ": " + format + "\n")
```

(b) (10 points) Output

> **Solution:**
>
> | Date | Format |
> |------------|-----------|
> | 05/12/2023 | ambiguous |
> | 12/05/2023 | ambiguous |
> | 03/08/2024 | ambiguous |
> | 08/03/2024 | ambiguous |
> | 05/06/2023 | ambiguous |

(c) (20 points) Report

(a) Proposed Algorithm

> **Solution:** The proposed algorithm extracts dates from the corpus using a regular expression. It then classifies the dates based on the following criteria:
>
> - If the first two digits are greater than 12, then the date format is DD/MM/YYYY.
>
> - If the second pair of digits is greater than 12, then the date format is MM/DD/YYYY.
>
> - If both the pairs are less than 12, then the code checks for contextual clues to understand the date format i.e. if the corpus uses US words, then the format is MM/DD/YYYY and if it uses British words, then the format is DD/MM/YYYY.
>
> - In case the dates are not classified by any of the above conditions, the dates are labelled ambiguous.

(b) List of Extracted Dates and Interpretation

> **Solution:**
>
> | Date | Interpretation |
> |------------|----------------|
> | 05/12/2023 | ambiguous |
> | 12/05/2023 | ambiguous |
> | 03/08/2024 | ambiguous |
> | 08/03/2024 | ambiguous |
> | 05/06/2023 | ambiguous |

(c) Challenges

> **Solution:** Logical clues were not enough to classify the dates in the given corpus. Therefore, I had to rely on the contextual clues. For this, I generated a list of words that could be used to differentiate between a US English text and a British English text. However, I still ended up with date format being ambiguous. This could be because the used list of words is finite and might not have the words present within the corpus. Moreover, it is possible that the corpus does not have any such words that are different in US and English British. Hence, the formats remain ambiguous.

2. (40 points) Identifying the First 10 Merges in a Wordpiece Algorithm

   (a) (20 points) Python Script

   **Solution:**

```
1
2 import re
3
4 # --------------------------- Input Handling --------------------------- #
5 def input_handling(file_name):
6   corpus = open(file_name, "r")
7   corpus = corpus.read()
8   corpus = corpus.lower()
9   corpus = corpus.split()
10   return corpus
11
12 # ----------- Tokenization based on words + Frequency of each word -------- #
13 def calculate_word_freq(corpus):
14   word_freq = {}
15   for w in corpus:
16       if w in word_freq:
17           word_freq[w] += 1
18       else:
19           word_freq[w] = 1
20   return word_freq
21
22 # ------- Tokenization based on letters + Frequency of each letter ------- #
23 def letter_tokenization(word_freq):
24   tokenized = []
25   letter_freq = {}
26   for i in word_freq:
27       s = [i[0]]
28       if i[0] in letter_freq:
29           letter_freq[i[0]] += word_freq[i]
30       else:
31           letter_freq[i[0]] = word_freq[i]
32
33       for k in i[1:]:
34           s.append('##'+k)
35
36           if ('##'+k) in letter_freq:
37               letter_freq['##'+k] += word_freq[i]
38           else:
39               letter_freq['##'+k] = word_freq[i]
40
41       tokenized.append((s, word_freq[i]))
42   return ((tokenized, letter_freq))
43
44 # --------------------------- Merging ----------------------------------- #
45 def merge(tokenized, letter_freq, merge_limit):
46   merges = 0
47   vocab = []
48
49   while merges != merge_limit:
50       pairs_freq = {}
51
52       # --------- Finding pairs + Frequency of each pair ------------ #
53       for t in range(len(tokenized)):
54           for u in range(len(tokenized[t][0])-1):
55               new_pair = tokenized[t][0][u] + tokenized[t][0][u+1]
56               if new_pair not in pairs_freq:
57                   #[freq_pair, freq_1, freq_2, [(row, col)]]
58                   pairs_freq[new_pair] = [tokenized[t][1], letter_freq[
   tokenized[t][0][u]], letter_freq[tokenized[t][0][u+1]], [(t,u)]]
59               else:
60                   pairs_freq[new_pair][0] += tokenized[t][1]
61                   pairs_freq[new_pair][3] += [(t,u)]
62
63       # --------- Calculating score for each pair ---------------- #
64       score_arr = []
65       for p in pairs_freq:
66           score = pairs_freq[p][0] / (pairs_freq[p][1] * pairs_freq[p][2])
```

```
67            score_arr.append((p, score))
68
69        score_arr.sort(key=lambda x: x[1], reverse=True)
70
71
72        # ----------------- Updating data ------------------------ #
73
74        vocab.append(score_arr[0][0])
75        letter_freq[score_arr[0][0]] = pairs_freq[score_arr[0][0]][0]
76
77        change = pairs_freq[score_arr[0][0]][3]
78        for c in change:
79            tokenized[c[0]][0][c[1]] = tokenized[c[0]][0][c[1]] + tokenized[c
    [0]][0][c[1]+1]
80            tokenized[c[0]][0].remove(tokenized[c[0]][0][c[1]+1])
81
82        merges += 1
83    return vocab
84
85#------------------------- Main Function --------------------------------#
86
87def main():
88    file_name = "wordpiece_input.txt"
89    corpus = input_handling(file_name)
90    word_frequency = calculate_word_freq(corpus)
91    tokenized, letter_freq = letter_tokenization(word_frequency)
92    vocab = merge(tokenized, letter_freq, 10)
93
94    output_file = open("q2_batool_ba07612.txt", "w", encoding='utf-8')
95
96    output_file.write("Merged Vocabulary:\n\n")
97
98    for i in range(len(vocab)):
99        if vocab[i][0] == '#':
100            vocab[i] = vocab[i][0:2] + re.sub('#', '', vocab[i][2:])
101        else:
102            vocab[i] = re.sub('#', '', vocab[i])
103        output_file.write(vocab[i]+'\n')
104
105 main()
106
```

(b) (20 points) Report

(a) First 10 Merges

**Solution:**

| Merge Number | Token 1 | Token 2 | Merged Tokens |
|---|---|---|---|
| 1. | 1 | ##0 | 10 |
| 2. | e | ##x | ex |
| 3. | o | ##f | of |
| 4. | ##f | ##y | ##fy |
| 5. | ##m | ##p | ##mp |
| 6. | ##q | ##u | ##qu |
| 7. | ##b | ##u | ##bu |
| 8. | ##" | ##, | ##", |
| 9. | ex | ##a | exa |
| 10. | exa | ##mp | examp |

(b) Challenges

> **Solution:** Initially, I used a set structure to store my vocabulary. This caused the output to be different in every run because the order of the pairs kept changing. Moreover, multiple pairs had the same score. This also caused the output to differ in different runs. Therefore, I had to switch to a list structure to avoid this.

3. (20 points) Tokenizing Urdu Text

   (a) (5 points) Python Script

   **Solution:**

```python
corpus = open("urdu_text_input.txt", "r", encoding='utf-8')
corpus = corpus.read()

#--------------------- Using Regular Expressions ------------------------#
import re

word_tokens_re = set(re.findall(r'\w+', corpus)) #word tokenization
letter_tokens_re = set(re.findall(r'\w', corpus)) #letter tokenization

output_file = open("q3_batool_ba07612.txt", "w", encoding='utf-8')

output_file.write("Using Regular Expression\n\n")
word_tokens_str = "\n".join(word_tokens_re)
output_file.write("Word Based\n")
output_file.write(word_tokens_str)
letter_tokens_str = "\n".join(letter_tokens_re)
output_file.write("\n\nLetter Based\n")
output_file.write(letter_tokens_str)
output_file.write("\n\n\n")

#------------------------- Using nltk Library --------------------------#
from nltk.tokenize import word_tokenize
word_tokens_nltk = set(word_tokenize(corpus))
output_file.write("Using nltk Library\n")
output_file.write("\n".join(word_tokens_nltk))
output_file.write("\n\n\n")

#------------------------ Using split() Method -------------------------#

word_tokens_split = set(corpus.split())
output_file.write("Using split() Method\n")
output_file.write("\n".join(word_tokens_split))
```

(b) (10 points) Output

**Solution:**

- Using Regular Expression
  - Word Based Tokenization



  - Character Based Tokenization



- Using nltk Library



- Using split() Method



(c) (10 points) Report
   (a) Code

**Solution:**

```python
corpus = open("urdu_text_input.txt", "r", encoding='utf-8')
corpus = corpus.read()

#---------------------- Using Regular Expressions --------------------#
import re

word_tokens_re = set(re.findall(r'\w+', corpus)) #word tokenization
letter_tokens_re = set(re.findall(r'\w', corpus)) #letter tokenization

output_file = open("q3_batool_ba07612.txt", "w", encoding='utf-8')

output_file.write("Using Regular Expression\n\n")
word_tokens_str = "\n".join(word_tokens_re)
output_file.write("Word Based\n")
output_file.write(word_tokens_str)
letter_tokens_str = "\n".join(letter_tokens_re)
output_file.write("\n\nLetter Based\n")
output_file.write(letter_tokens_str)
output_file.write("\n\n\n")

#---------------------- Using nltk Library ----------------------#
from nltk.tokenize import word_tokenize
```

```
24 word_tokens_nltk = set(word_tokenize(corpus))
25 output_file.write("Using nltk Library\n")
26 output_file.write("\n".join(word_tokens_nltk))
27 output_file.write("\n\n\n")
28
29 #------------------------ Using split() Method ----------------------#
30
31 word_tokens_split = set(corpus.split())
32 output_file.write("Using split() Method\n")
33 output_file.write("\n".join(word_tokens_split))
34
```

(b) Challenges

**Solution:** I tried using different methods for tokenizing the text. The nltk library and the built in split() method produced the same output. The output from the Regular Expressions was a little different. The word tokenization using Regular Expressions ignored the punctuation marks while nltk and split() methods included the punctuation marks and treated them as a part of the word before them. I also tried to tokenize the text using the UrduHack library but the library had some dependencies on other libraries which were not available.