

CSE517 Assignment: A2_1.3

A well-known problem with substitution ciphers is that, in a long ciphertext, the frequency of $\text{encrypt}(x)$ will be close to the frequency of x in plaintext. In natural languages, there is a lot of variance in different symbols' frequencies. (For words, "Zipf's Law" states that the probability of the r th most frequent word in a text corpus will have relative frequency proportional to $1/r$.) Your task is to exploit this problem to decrypt the ciphertext we provide to you. The original plaintext is in all-upper-cased English. You may assume the space symbol, numerical digits, and any other non-alphabetic symbols are fixed. This implies that $N = 26$. We suggest that you automate the calculation and visualization of the symbol frequencies. Some manual search may be required since your ciphertext (and any plaintext you use to estimate English relative frequencies) will be finite, leading to variance in your estimates. You should submit your decrypted text.

```
In [38]: 1 import random
          2 import numpy as np
          3 import matplotlib.pyplot as plt
          4 import pandas as pd
          5 from math import sqrt
```

```
In [11]: 1 e=open(r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A2/A2.encrypted')
```

```
In [12]: 1 def convert_line_to_dict(input_line, char_dict):
          2     for char in input_line:
          3         if char in char_dict:
          4             char_dict[char] = char_dict.get(char) + 1
          5         else:
          6             if (ord(char) >= 65) and (ord(char) <= 90):
          7                 char_dict[char] = 1
          8
          9 def convert_text_to_dict(file, char_dict):
          10     for line in file:
          11         convert_line_to_dict(line, char_dict)
```

```
In [13]: 1 encrypted_char_dict = {}
          2 convert_text_to_dict(e, encrypted_char_dict)
          3 encrypted_char_list = sorted(encrypted_char_dict.items(), key=lambda x:
          4 encrypted_char_list
```

```
Out[13]: [('O', 495),
           ('H', 357),
           ('Z', 313),
           ('L', 288),
           ('B', 283),
           ('F', 277),
           ('T', 266),
           ('X', 209),
           ('D', 182),
           ('V', 170),
           ('C', 159),
           ('Y', 126),
           ('I', 111),
           ('R', 97),
           ('N', 88),
           ('P', 84),
           ('U', 79),
           ('E', 76),
           ('G', 76),
           ('S', 73),
           ('K', 46),
           ('W', 42),
           ('M', 9),
           ('Q', 6),
           ('A', 5),
           ('J', 4)]
```

Train Corpus: <https://www.gutenberg.org/files/65728/65728-0.txt>
(<https://www.gutenberg.org/files/65728/65728-0.txt>)

```
In [14]: 1 d=open(r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A2/A2_train_cor
```

```
In [15]: 1 plain_char_dict = {}  
2 convert_text_to_dict(d, plain_char_dict)  
3 plain_char_list = sorted(plain_char_dict.items(), key=lambda x: x[1], r  
4 plain_char_list
```

```
Out[15]: [('E', 93834),  
          ('T', 75541),  
          ('O', 66465),  
          ('I', 56202),  
          ('R', 56152),  
          ('N', 52641),  
          ('A', 52318),  
          ('S', 42572),  
          ('H', 33268),  
          ('D', 31694),  
          ('L', 27999),  
          ('U', 25138),  
          ('C', 23638),  
          ('G', 19864),  
          ('M', 18199),  
          ('P', 18062),  
          ('Y', 17649),  
          ('F', 16902),  
          ('W', 14910),  
          ('B', 12856),  
          ('K', 7796),  
          ('V', 5374),  
          ('J', 3106),  
          ('Q', 1586),  
          ('X', 1242),  
          ('Z', 62)]
```

```

In [16]: 1 converter_dict = {}
          2 for i in range(26):
          3     converter_dict[encrypted_char_list[i][0]] = plain_char_list[i][0]
          4 # converter_dict['I'] = 'S'
          5 converter_dict['F'] = 'S'
          6 converter_dict['E'] = 'B'
          7 converter_dict['N'] = 'U'
          8 converter_dict['C'] = 'U'
          9 converter_dict['B'] = 'A'
         10 converter_dict['U'] = 'G'
         11 converter_dict['T'] = 'N'
         12 converter_dict['Y'] = 'M'
         13 converter_dict['S'] = 'P'
         14 converter_dict['I'] = 'D'
         15 converter_dict['P'] = 'F'
         16 converter_dict['K'] = 'V'
         17 converter_dict['V'] = 'H'
         18 converter_dict['X'] = 'R'
         19 converter_dict['G'] = 'Y'
         20 converter_dict['R'] = 'C'
         21 converter_dict['D'] = 'L'
         22 converter_dict['Q'] = 'Z'
         23 converter_dict['W'] = 'K'
         24 converter_dict['J'] = 'Q'
         25 converter_dict['N'] = 'W'
         26 # converter_dict['V'] = 'H'
         27 converter_dict

```

```

Out[16]: {'O': 'E',
          'H': 'T',
          'Z': 'O',
          'L': 'I',
          'B': 'A',
          'F': 'S',
          'T': 'N',
          'X': 'R',
          'D': 'L',
          'V': 'H',
          'C': 'U',
          'Y': 'M',
          'I': 'D',
          'R': 'C',
          'N': 'W',
          'P': 'F',
          'U': 'G',
          'E': 'B',
          'G': 'Y',
          'S': 'P',
          'K': 'V',
          'W': 'K',
          'M': 'J',
          'Q': 'Z',
          'A': 'X',
          'J': 'Q'}

```

```
In [17]: 1 e=open(r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A2/A2.encrypted')
2 o=open(r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A2/A2.decrypted')
3
4 def decrypt_line(input_line):
5     output_line = ""
6     for char in input_line:
7         if (ord(char) >= 65) and (ord(char) <= 90):
8             output_line = output_line + converter_dict[char]
9         else:
10            output_line = output_line + char
11    return output_line
12
13 def decrypt_text(file):
14     for line in file:
15         nline = decrypt_line(line)
16         o.write(nline)
17    o.close()
18
19 decrypt_text(e)
```

```
In [21]: 1 total_count = 0
2 encrypted_char_list_with_freq = []
3 for element in encrypted_char_list:
4     total_count = total_count + element[1]
5 for element in encrypted_char_list:
6     encrypted_char_list_with_freq.append((converter_dict[element[0]], e
7 encrypted_char_list_with_freq
```

```
Out[21]: [('E', 0.1262433052792655),
('T', 0.09104820198928845),
('O', 0.07982657485335373),
('I', 0.07345065034429993),
('A', 0.07217546544248916),
('S', 0.07064524356031625),
('N', 0.06783983677633257),
('R', 0.05330272889568988),
('L', 0.04641673042591176),
('H', 0.04335628666156593),
('U', 0.04055087987758225),
('M', 0.032134659525631215),
('D', 0.028309104820198928),
('C', 0.024738587095128793),
('W', 0.022443254271869422),
('F', 0.02142310635042081),
('G', 0.020147921448610048),
('B', 0.01938281050752359),
('Y', 0.01938281050752359),
('P', 0.018617699566437135),
('V', 0.011731701096659015),
('K', 0.010711553175210406),
('J', 0.0022953328232593728),
('Z', 0.001530221882172915),
('X', 0.0012751849018107625),
('Q', 0.00102014792144861)]
```

```
In [34]: 1 encrypter_dict = {v: k for k, v in converter_dict.items()}
          2 encrypter_dict
```

```
Out[34]: {'E': 'O',
          'T': 'H',
          'O': 'Z',
          'I': 'L',
          'A': 'B',
          'S': 'F',
          'N': 'T',
          'R': 'X',
          'L': 'D',
          'H': 'V',
          'U': 'C',
          'M': 'Y',
          'D': 'I',
          'C': 'R',
          'W': 'N',
          'F': 'P',
          'G': 'U',
          'B': 'E',
          'Y': 'G',
          'P': 'S',
          'V': 'K',
          'K': 'W',
          'J': 'M',
          'Z': 'Q',
          'X': 'A',
          'Q': 'J'}
```

Part 1.5: Lets add more symbols '@', '#', '\$', '%', '~', '^'

$|EUT| = 26 + 6 \text{ (new ones)} = 32$

Most frequent symbols: Number of symbols for 'E' = 3 Number of symbols for 'T' = 2 Number of symbols for 'O' = 2 Number of symbols for 'I' = 2 Number of symbols for 'A' = 2

```
In [35]: 1 encrypter_dict['E'] = ('O', '@', '#')
          2 encrypter_dict['T'] = ('H', '%')
          3 encrypter_dict['O'] = ('Z', '~')
          4 encrypter_dict['I'] = ('L', '^')
          5 encrypter_dict['A'] = ('B', '$')
```

In [36]: 1 encrypter_dict

```
Out[36]: {'E': ('O', '@', '#'),
          'T': ('H', '%'),
          'O': ('Z', '~'),
          'I': ('L', '^'),
          'A': ('B', '$'),
          'S': 'F',
          'N': 'T',
          'R': 'X',
          'L': 'D',
          'H': 'V',
          'U': 'C',
          'M': 'Y',
          'D': 'I',
          'C': 'R',
          'W': 'N',
          'F': 'P',
          'G': 'U',
          'B': 'E',
          'Y': 'G',
          'P': 'S',
          'V': 'K',
          'K': 'W',
          'J': 'M',
          'Z': 'Q',
          'X': 'A',
          'Q': 'J'}
```

```
In [45]: 1 e=open(r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A2/A2.decrypted
2 o=open(r'/Users/nehakardam/Documents/UWclasses /CSE NLP/A2/A2.new_encry
3 def encrypt_line(input_line):
4     output_line = ""
5     for char in input_line:
6         if (ord(char) >= 65) and (ord(char) <= 90):
7             if type(encrypter_dict['E']) is tuple:
8                 output_line = output_line + random.choice(encrypter_dic
9             else:
10                output_line = output_line + encrypter_dict[char]
11        else:
12            output_line = output_line + char
13    return output_line
14
15 def encrypt_text(file):
16     for line in file:
17         nline = encrypt_line(line)
18         o.write(nline)
19     o.close()
20
21 encrypt_text(e)
```

