

1.9 Hangman

```
In [1]: import pandas as pd
import string
```

```
In [2]: # Read in word counts file
word_counts = pd.read_csv('./hw1_word_counts_05.txt', sep=' ',
                           header=None, names=['word', 'count'])

# Compute probabilities of words
word_counts['P(W=w)'] = word_counts['count'] / word_counts['count'].sum(
axis=0)

word_counts.head()
```

Out[2]:

	word	count	P(W=w)
0	AARON	413	0.000054
1	ABABA	199	0.000026
2	ABACK	64	0.000008
3	ABATE	69	0.000009
4	ABBAS	290	0.000038

(a) Print out 15 most frequent and 14 least frequent 5-letter words.

Compute the prior probabilities $P(w) = COUNT(w)/COUNT_{total}$. As a sanity check, print out the 15 most frequent 5-letter words, as well as the 14 least frequent words.

```
In [3]: top_15 = word_counts.sort_values(by='P(W=w)', ascending=False)
top_15.head(15)
```

Out[3]:

	word	count	P(W=w)
5821	THREE	273077	0.035627
5102	SEVEN	178842	0.023333
1684	EIGHT	165764	0.021626
6403	WOULD	159875	0.020858
18	ABOUT	157448	0.020542
5804	THEIR	145434	0.018974
6320	WHICH	142146	0.018545
73	AFTER	110102	0.014365
1975	FIRST	109957	0.014346
1947	FIFTY	106869	0.013943
4158	OTHER	106052	0.013836
2073	FORTY	94951	0.012388
6457	YEARS	88900	0.011598
5806	THERE	86502	0.011286
5250	SIXTY	73086	0.009535

```
In [4]: last_14 = word_counts.sort_values(by='P(W=w)')
last_14.head(14)
```

Out[4]:

	word	count	P(W=w)
3554	MAPCO	6	7.827935e-07
712	BOSAK	6	7.827935e-07
895	CAIXA	6	7.827935e-07
4160	OTTIS	6	7.827935e-07
5985	TROUP	6	7.827935e-07
1107	CLEFT	7	9.132590e-07
2041	FOAMY	7	9.132590e-07
977	CCAIR	7	9.132590e-07
5093	SERNA	7	9.132590e-07
6443	YALOM	7	9.132590e-07
5872	TOCOR	7	9.132590e-07
3978	NIAID	7	9.132590e-07
4266	PAXON	7	9.132590e-07
1842	FABRI	7	9.132590e-07

Do your results make sense?

Yes. The most common words are ones that are often used including numbers. The least common words are very specific and not used in daily speech.

(b) The Best Next Guess

Consider the following stages of the game. For each of the following, indicated the best next guess -- namely, the letter l that is most probable to be among the missing letters. Also report the probability $P(L_i = l \text{ for some } i \text{ in } 1, 2, 3, 4, 5 | E)$ for your guess l . Your answers should fill in the last two columns of this table.

```

In [22]: def p_e_given_w(correct , incorrect):
'''
    Function to give  $P(E|W)$  for each word. Checks if the guessed correctly/incorrectly
    matches with each word.

    correct: list of correct guesses
    incorrect: list of incorrect guesses

    p_e_w: list of probabilities  $P(E|W)$  for each word
'''
# Initialize  $P(E|W)$  for each word
p_e_w = [None]*len(word_counts.index)

# Compare words to correct/incorrect guesses to determine  $P(E|W)$ 
for i in range(len(word_counts.index)):
    for letter, guess in zip(word_counts['word'][i], correct):
        #  $P(E|W) = 0$  if the word contains a letter that is "incorrect"
        if letter in incorrect:
            p_e_w[i] = 0
            break

        #  $P(E|W) = 0$  if the letter is in the guessed spot but the
        # word doesn't have the same letter in that spot.
        elif (guess != None) and (letter != guess):
            p_e_w[i] = 0
            break

        #  $P(E|W) = 0$  if a letter has been guessed and is in the word
        # but not in the same positions as the word in question.
        elif (letter in correct) and (guess == None):
            p_e_w[i] = 0
            break

    else:
        p_e_w[i] = 1

return p_e_w

```

```
In [23]: def p_l_given_w(letter):  
        '''  
        Function that gives the probabilities of a certain letter given a word.  
  
        letter: letter of interest from alphabet  
  
        p_l_w: probability of that letter for each word  
        '''  
        # Initialize P(L/W) for each word  
        p_l_w = [None]*len(word_counts.index)  
  
        # Check if letter is anywhere in word  
        for i in range(len(word_counts.index)):  
            if letter in list(word_counts['word'][i]):  
                p_l_w[i] = 1  
            else:  
                p_l_w[i] = 0  
  
        return p_l_w
```

```

In [31]: def best_guess(correct, incorrect):
    '''
    Gives the best next guess l and P(L=l|E)

    correct: list of correct guesses in order
    incorrect: list of incorrect guesses
    '''

    # Compute P(E|W=w)
    word_counts['P(E|W=w)'] = p_e_given_w(correct, incorrect)

    # Compute P(W=w|E)
    word_counts['P(E|W=w)*P(W=w)'] = word_counts['P(E|W=w)'] * word_counts['P(W=w)']
    word_counts['P(W=w|E)'] = word_counts['P(E|W=w)*P(W=w)'] / word_counts['P(E|W=w)*P(W=w)'].sum(axis=0)
    del word_counts['P(E|W=w)*P(W=w)']

    # Compute P(L=l|W=w)
    for letter in string.ascii_uppercase:
        word_counts['P(L={}|W=w)'.format(letter)] = p_l_given_w(letter)

    # Compute P(L=l|W=w)*P(W=w|E) for each letter and word
    for letter in string.ascii_uppercase:
        word_counts['P(L={}|W=w)*P(W=w|E)'.format(letter)] = p_l_given_w(letter) * word_counts['P(W=w|E)']

    # Compute P(L=l|E)
    p_l_e = [None]*len(string.ascii_lowercase)
    for letter, i in zip(string.ascii_uppercase, range(len(string.ascii_lowercase))):
        p_l_e[i] = word_counts['P(L={}|W=w)*P(W=w|E)'.format(letter)].sum(axis=0)

    p_letter_e = pd.DataFrame(p_l_e, columns=['P(L=l|E)'], index=list(string.ascii_uppercase))

    best_guess = p_letter_e.loc[p_letter_e['P(L=l|E)'] < 0.99999].idxmax().index.values[0]
    max_p_l_e = p_letter_e.loc[p_letter_e['P(L=l|E)'] < 0.99999].idxmax().values[0][0]

    print('For correct guesses', correct, 'and incorrect guesses {}'.format(incorrect))
    print('Your best next guess is', best_guess, 'with a probability P(L={}|E) of'.format(best_guess), round(max_p_l_e, 4), '\n')

```

```
In [32]: # Check against given solutions
correct = [None] * 5
incorrect = ['E', 'O']
best_guess(correct, incorrect)

correct = ['D', None, None, 'I', None]
incorrect = []
best_guess(correct, incorrect)

incorrect = ['A']
best_guess(correct, incorrect)

correct = [None, 'U', None, None, None]
incorrect = ['A', 'E', 'I', 'O', 'S']
best_guess(correct, incorrect)
```

For correct guesses [None, None, None, None, None] and incorrect guesses ['E', 'O']:

Your best next guess is I with a probability $P(L=I|E)$ of 0.6366

For correct guesses ['D', None, None, 'I', None] and incorrect guesses []:

Your best next guess is A with a probability $P(L=A|E)$ of 0.8207

For correct guesses ['D', None, None, 'I', None] and incorrect guesses ['A']:

Your best next guess is E with a probability $P(L=E|E)$ of 0.7521

For correct guesses [None, 'U', None, None, None] and incorrect guesses ['A', 'E', 'I', 'O', 'S']:

Your best next guess is Y with a probability $P(L=Y|E)$ of 0.627

```
In [33]: # New solutions
correct = [None] * 5
incorrect = []
best_guess(correct, incorrect)

incorrect = ['A', 'I']
best_guess(correct, incorrect)

correct = ['A', None, None, None, 'R']
incorrect = []
best_guess(correct, incorrect)

incorrect = ['E']
best_guess(correct, incorrect)

correct = [None, None, 'U', None, None]
incorrect = ['O', 'D', 'L', 'C']
best_guess(correct, incorrect)
```

For correct guesses [None, None, None, None, None] and incorrect guesses []:

Your best next guess is E with a probability $P(L=E|E)$ of 0.5394

For correct guesses [None, None, None, None, None] and incorrect guesses ['A', 'I']:

Your best next guess is E with a probability $P(L=E|E)$ of 0.6214

For correct guesses ['A', None, None, None, 'R'] and incorrect guesses []:

Your best next guess is T with a probability $P(L=T|E)$ of 0.9816

For correct guesses ['A', None, None, None, 'R'] and incorrect guesses ['E']:

Your best next guess is O with a probability $P(L=O|E)$ of 0.9913

For correct guesses [None, None, 'U', None, None] and incorrect guesses ['O', 'D', 'L', 'C']:

Your best next guess is T with a probability $P(L=T|E)$ of 0.7045

In []: