

Natural Language Processing

Lab week 3 – Language modelling



The objective of these practical exercises is to familiarize yourself with NLP programming using python.

1. Frequency

Given an input sentence, calculate the **frequency** ($p(w)$) of each word (w) in the sentence according to the formula:

$$p(w) = \frac{\text{occurrences of word}}{\text{number of tokens}} \quad (1)$$

Input: “the cat sat on the mat with a cat”

Output:

The word “a” frequency is: 0.1111111111
The word “on” frequency is: 0.1111111111
The word “mat” frequency is: 0.1111111111
The word “cat” frequency is: 0.2222222222
The word “the” frequency is: 0.2222222222
The word “with” frequency is: 0.1111111111
The word “sat” frequency is: 0.1111111111

2. Unigram LM

Given an input sentence (s), calculate the **unigram language model** of the sentence according to the formula:

$$p(s = w_1, \dots, w_n) = p(w_1) \times \dots \times p(w_n) \quad (2)$$

Hint : Interpolation of the $P(w)$ function in Question 1 could be a good idea.

Input: “the cat sat on the mat with a cat”

Output: 8.36300632515e-07

3. Bigram LM

Following Question 1&2, write a program to compute **bigram probability of a sentence**. The input to your program is a file containing a number of sentences and the output is the probability of one sentence. To compute **bigram relative frequency** use this formula:

$$p(w_2|w_1) = \frac{\text{count}(w_1, w_2)}{\sum_w \text{count}(w_1, w)} \quad (3)$$

To compute the bigram probability of a sentence use this formula:

$$p(s) = p(w_2|w_1) \times p(w_3|w_2) \dots \times p(w_n|w_{n-1}) \quad (4)$$

Hint:

1, Interpolation of the function in Question 1 of Lab-3 could be a good idea. 2, Creating functions based on Question 1 and 2 could be a good idea.

Input: file_name.txt

Calculate the probability of the sentence "<s> a cat sat on the mat </s>"

Output: 0.00097615576843

4. Smoothing

First, try another sentence using your program of Question 3:

Calculate the probability of the sentence "<s> a cat sat on the car </s>". What result do you get?

Think about what the reason is and why we need smoothing technique in language modelling.

Second, modify your function of **bigram relative frequency** according to add-one smoothing formula:

$$p(w_2|w_1) = \frac{\text{count}(w_1, w_2) + 1}{\sum_w \text{count}(w_1, w) + v} \quad (5)$$

where v is vocabulary size (how many unique words in your file). Use your smoothed function to calculate **bigram probability of a sentence** of the two sentences.

Input: file_name.txt

Calculate the probability of the sentence "<s> a cat sat on the mat </s>"

Output: 0.000140949604457

Calculate the probability of the sentence "<s> a cat sat on the car </s>"

Output: 3.00170453936e-05

Optional-

In order to adapt your **bigram probability** program to **n-gram probability** program.

Add one more input to your program of Question 4.

Input:

1, file_name.txt 2, gram_number

Calculate the n-gram probability of the sentence "< s> a cat sat on the mat </s>".

Output:

1-gram: 2.28175851587e-08

2-gram: 0.000140949604457

3-gram: 0.000263061746438

4-gram: 0.000423106305459