

# HW 11 보고서

2013-10035  
언어학과 김영민

1.1

# 1.1 문제 해결 방법

```
import numpy as np
# corpus size
N = 1410000000

# N-gram frequency
ngram_dict = {}
ngram_dict['1'] = {'하': 3520000,
                  '파': 392000,
                  '단': 34600,
                  '에': 339000,
                  '로': 24300,
                  '고': 359000}

ngram_dict['2'] = {'하': 56100,
                  '파': 23,
                  '단': 160,
                  '에': 85,
                  '로': 198}

ngram_dict['3'] = {'하': 34,
                  '파': 0,
                  '단': 3,
                  '에': 85}

def ngram_prob(ngram):
    if '<s>' in ngram:
        return ngram_prob(ngram[1:])
    else:
        n = len(ngram)
        if n == 1:
            prob = ngram_dict[str(n)].get(ngram[0], 0) / N
            print('P(%s): %.32f' % (ngram[0], prob))
            return prob
        prob = ngram_dict[str(n)].get(' '.join(ngram), 0) / ngram_dict[str(n-1)].get(' '.join(ngram[:-1]), 0)
        print('P(%s): %.32f' % (' '.join(ngram), prob))
        return prob

def sentence_probability(n, sentence):
    sentence = ['<s>'] * (n-1) + sentence.split(' ')
    prob = 1
    for i in range(0, len(sentence) - n + 1):
        ngram = sentence[i: i + n]
        prob *= ngram_prob(ngram)
    return prob
```

ngram.py 파일에 이와 같은 메소드를 작성해서  $n(=1, 2, 3)$  을 파라미터로 줬을 때 각각 unigram, bigram, trigram 모델 기준으로 문장의 확률을 계산할 수 있도록 함.

## 1.1 메소드 설명 sentence\_probability()

```
def sentence_probability(n, sentence):  
    sentence = ['<s>'] * (n-1) + sentence.split(' ')  
    prob = 1  
    for i in range(0, len(sentence) - n + 1):  
        ngram = sentence[i: i + n]  
        prob *= ngram_prob(ngram)  
    return prob
```

- 문장을 공백문자로 split 해서 list 로 만들고 n값에 따라 적절한 개수의 <s> 문자를 list 의 앞에 삽입
- 전체 문장을 n-gram 단위로 잘라서 각 n-gram 의 확률(ngram\_prob())을 구한 뒤 곱해줌
- 출력: 전체 문장 W 의 확률  $P(W)$

## 1.1 메소드 설명 ngram\_prob()

```
def ngram_prob(ngram):
    if '<s>' in ngram:
        return ngram_prob(ngram[1:])
    else:
        n = len(ngram)
        if n == 1:
            prob = ngram_dict[str(n)].get(ngram[0], 0) / N
            print('P(%s): %.32f' % (ngram[0], prob))
            return prob
        prob = ngram_dict[str(n)].get(' '.join(ngram), 0) / ngram_dict[str(n-1)].get(' '.join(ngram[:-1]), 0)
        print('P(%s): %.32f' % (' '.join(ngram), prob))
        return prob
```

- 문장을 n-gram 단위로 나눴을 때 각 n-gram 의 확률을 구하는 메소드
- 문제에서 <s> 가 포함된 n-gram 은 <s>를 모두 제거한 (n-k)-gram 과 확률이 같다는 가정이 있으므로 <s> 가 있으면 제거하고 확률을 리턴하도록 재귀적으로 처리

## 1.1 결과

```
P(하늘은) : 0.00249645390070921992242691800357
P(파랑고) : 0.00027801418439716310505513408025
P(단풍잎은) : 0.00002453900709219858001131249481
P(빨강고) : 0.000240425531914893617279717144645
P(행복은) : 0.00001723404255319148806820279962
P(노랑고) : 0.00025460992907801416286137086153

>> UNIGRAM:
PROBABILITY: 1.7967620814451954e-23
PERPLEXITY: 6179.0006046138515
=====
P(하늘은) : 0.00249645390070921992242691800357
P(파랑고) : 0.015937500000000000027755575615629
P(단풍잎은) : 0.0000586734693877551035387854056
P(빨강고) : 0.00462427745664739913572516272211
P(행복은) : 0.00025073746312684363265338438609
P(노랑고) : 0.00814814814814814741827930788531

>> BIGRAM:
PROBABILITY: 2.2055024554912734e-17
PERPLEXITY: 597.1481287638287
=====
P(하늘은) : 0.00249645390070921992242691800357
P(하늘고) : 0.015937500000000000027755575615629
P(하늘단풍잎은) : 0.00000606060606060606008038682546
P(하늘단풍고) : 0.0000000000000000000000000000000000
P(단풍잎은) : 0.0187499999999999999999999999999999
P(단풍고) : 1.0000000000000000000000000000000000
P(단풍빨강고) : 0.0000000000000000000000000000000000
P(단풍행복은) : 0.0000000000000000000000000000000000
P(단풍노랑고) : 0.0000000000000000000000000000000000

>> TRIGRAM:
PROBABILITY: 0.0
Traceback (most recent call last):
  File "ngram.py", line 78, in <module>
    print('PERPLEXITY:', perplexity(trigram_prob, sent))
  File "ngram.py", line 49, in perplexity
    return probability ** (-1/k)
ZeroDivisionError: 0.0 cannot be raised to a negative power
```

- $n = 1, 2, 3$ ,  $\text{add\_k} = 1$ 로 코드를 돌린 결과

Unigram: 1.7967620814451954e-23

Bigram: 2.2055024554912734e-17

Trigram: 0.0

- ‘파랗고 단풍잎은 빨갛고’의 빈도가 0  
이므로  $P(W) = 0$ 의 결과가 나온다.

## 1.1 결과 검증 (Uni-gram)

```
>>> # UNIGRAM
...
>>> N = 1410000000
>>> p_1 = 3520000 / N # 하늘은
>>> p_1
0.00249645390070922
>>> p_2 = 392000 / N # 파랗고
>>> p_2
0.0002780141843971631
>>> p_3 = 34600 / N # 단풍잎은
>>> p_3
2.453900709219858e-05
>>> p_4 = 339000 / N # 빨갛고
>>> p_4
0.00024042553191489362
>>> p_5 = 24300 / N # 은행잎은
>>> p_5
1.7234042553191488e-05
>>> p_6 = 359000 / N # 노랗고
>>> p_6
0.00025460992907801416
>>> p_1 * p_2 * p_3 * p_4 * p_5 * p_6
1.7967620814451954e-23
```

코드가 의도한대로 돌아가는지  
확인해보기 위해 파이썬 shell 로 직접  
계산 결과를 구한 뒤 비교해본 결과  
메소드를 사용해 구한 값과 일치하는  
것을 확인할 수 있었다.

## 1.1 결과 검증 (Bi-gram)

```
>>> # BIGRAM 수정 보기 삽입 서식 슬라이드 정렬
...
>>> p_1 = 3520000 / N # <s> 하늘은 == 하늘은
>>> p_1
0.00249645390070922
>>> p_2 = 56100 / 3520000 # 하늘은 파랗고
>>> p_2
0.0159375
>>> p_3 = 23 / 392000 # 파랗고 단풍잎은
>>> p_3
5.8673469387755104e-05
>>> p_4 = 160 / 34600 # 단풍잎은 빨강고
>>> p_4
0.004624277456647399
>>> p_5 = 85 / 339000 # 빨강고 은행잎은
>>> p_5
0.00025073746312684363
>>> p_6 = 198 / 24300 # 은행잎은 노랗고
>>> p_6
0.008148148148148147
>>> p_1 * p_2 * p_3 * p_4 * p_5 * p_6
2.2055024554912734e-17
```



## 1.1 결과 확인 (Tri-gram)

```
>>> # TRIGRAM
...
>>> p_1 = 3520000 / N # <s> <s> 하늘은 == <s> 하늘은 == 하늘은
>>> p_1
0.00249645390070922
>>> p_2 = 56100 / 3520000 # <s> 하늘은 파랗고 == 하늘은 파랗고
>>> p_2
0.0159375
>>> p_3 = 34 / 56100 # 하늘은 파랗고 단풍잎은
>>> p_3
0.0006060606060606061
>>> p_4 = 0 / 23 # 파랗고 단풍잎은 빨강고
>>> p_4
0.0
>>> p_5 = 3 / 160 # 단풍잎은 빨강고 은행잎은
>>> p_5
0.01875
>>> p_6 = 85 / 85 # 빨강고 은행잎은 노랗고
>>> p_6
1.0
>>> p_1 * p_2 * p_3 * p_4 * p_5 * p_6
0.0
```

1.2

## 1.2 메소드 정의

```
def perplexity(probability, sentence):  
    k = len(sentence.split(' '))  
    return probability ** (-1/k)
```

- 강의 슬라이드에 나온 복잡도 계산식을 그대로 코드로 옮긴 `perplexity()` 메소드를 작성하였다.
- 문장을 공백 기준으로 `split` 하여 문장의 길이 (단어수) `k`를 구한 뒤 각 모델을 통해 구한 문장 전체의 확률에  $-1/k$  승한다.

## 1.2 결과

```
=====
>> UNIGRAM:
    PROBABILITY: 1.7967620814451954e-23
    PERPLEXITY: 6179.0006046138515
=====

>> BIGRAM:
    PROBABILITY: 2.2055024554912734e-17
    PERPLEXITY: 597.1481287638287
=====

>> TRIGRAM:
    PROBABILITY: 0.0
Traceback (most recent call last):
  File "ngram.py", line 78, in <module>
    print('          PERPLEXITY:', perplexity(trigram_prob, sentence))
  File "ngram.py", line 49, in perplexity
    return probability ** (-1/k)
ZeroDivisionError: 0.0 cannot be raised to a negative power
```

- 1.1 에서 구한 확률을 perplexity() 함수에 넣어 얻은 복잡도는 이와 같다.
- unigram 보다 bigram 이 확률이 높기 때문에 bigram의 복잡도가 더 낮은 것을 확인할 수 있다.
- trigram은 확률이 0.0 이기 때문에  $0.0^{(-\frac{1}{k})}$  을 계산할 때 에러가 발생한다. (division by zero)

1.3

## 1.3 메소드 정의

```
def trigram_interpolation(sentence, weights):  
    print('weights: %s\n' % str(weights))  
    sentence = ['<s>'] * 2 + sentence.split(' ')  
    prob = 1  
    for i in range(0, len(sentence) - 2):  
        tri, bi, uni = sentence[i:i+3], sentence[i+1:i+3], sentence[i+2:i+3]  
        model_probs = ngram_prob(tri), ngram_prob(bi), ngram_prob(uni)  
        print('weighted probability of %s: %.32f\n' % (" ".join(tri), np.sum([p*w for p, w in zip(model_probs, weights)])))  
        prob *= np.sum([p*w for p, w in zip(model_probs, weights)])  
    return prob
```

- weights ( $w_{tri}$ ,  $w_{bi}$ ,  $w_{uni}$ ) 가 주어졌을 때 이 값들을 기준으로 보간법을 적용한 확률을 구하는 메소드를 이와 같이 구현하였다.
- 1.1 을 풀 때 사용했던 k-smoothing 기법은 사용하지 않았다. (k 값을 0 으로 설정)

# 1.3 결과 확인

weights: (0.5, 0.3, 0.2)

P(하늘은): 0.00249645390070921992242691800357  
P(하늘은 파라고): 0.00249645390070921992242691800357  
P(하늘은 파랑고): 0.00249645390070921992242691800357  
weighted probability of <s> <s> 하늘은: 0.00249645390070921992242691800357

P(하늘은 파라고): 0.01593750000000000027755575615629  
P(하늘은 파랑고): 0.01593750000000000027755575615629  
P(파랑고): 0.00027801418439716310505513408025  
weighted probability of <s> 하늘은 파라고: 0.01280560283687943307073808796304

P(하늘은 파라고 단풍잎은): 0.00060606060606060606008038682546  
P(파라고 단풍잎은): 0.00005867346938775510353878545056  
P(단풍잎은): 0.00002453900709219858001131249481  
weighted probability of 하늘은 파라고 단풍잎은: 0.00032554014526506926422919074859

P(파라고 단풍잎은 빨강고): 0.00000000000000000000000000000000  
P(단풍잎은 빨강고): 0.00462427745664739913572516272211  
P(빨강고): 0.00024042553191489361729717144645  
weighted probability of 파라고 단풍잎은 빨강고: 0.00143536834337719843707192879378

P(단풍잎은 빨강고 은행잎은): 0.01874999999999999930611060960928  
P(빨강고 은행잎은): 0.00025073746312684363265338438609  
P(은행잎은): 0.00001723404255319148806820279962  
weighted probability of 단풍잎은 빨강고 은행잎은: 0.00945366804744869095034065509253

P(빨강고 은행잎은 노랑고): 1.00000000000000000000000000000000  
P(은행잎은 노랑고): 0.00814814814814741827930788531  
P(노랑고): 0.00025460992907801416286137086153  
weighted probability of 빨강고 은행잎은 노랑고: 0.50249536643026004867351730354130

>> TRIGRAM INTERPOLATION:  
PROBABILITY: 7.096168276562839e-14  
PERPLEXITY: 155.41610326068232

weights: (0.7, 0.2, 0.1)

P(하늘은): 0.00249645390070921992242691800357  
P(하늘은 파라고): 0.00249645390070921992242691800357  
P(하늘은 파랑고): 0.00249645390070921992242691800357  
weighted probability of <s> <s> 하늘은: 0.00249645390070921948874604900936

P(하늘은 파라고): 0.01593750000000000027755575615629  
P(하늘은 파랑고): 0.01593750000000000027755575615629  
P(파랑고): 0.00027801418439716310505513408025  
weighted probability of <s> 하늘은 파라고: 0.01437155141843971493942344608286

P(하늘은 파라고 단풍잎은): 0.00060606060606060606008038682546  
P(파라고 단풍잎은): 0.00005867346938775510353878545056  
P(단풍잎은): 0.00002453900709219858001131249481  
weighted probability of 하늘은 파라고 단풍잎은: 0.00043843101882919507028199546106

P(파라고 단풍잎은 빨강고): 0.00000000000000000000000000000000  
P(단풍잎은 빨강고): 0.00462427745664739913572516272211  
P(빨강고): 0.00024042553191489361729717144645  
weighted probability of 파라고 단풍잎은 빨강고: 0.00094889804452096926205839633184

P(단풍잎은 빨강고 은행잎은): 0.01874999999999999930611060960928  
P(빨강고 은행잎은): 0.00025073746312684363265338438609  
P(은행잎은): 0.00001723404255319148806820279962  
weighted probability of 단풍잎은 빨강고 은행잎은: 0.01317687089688068628745565291638

P(빨강고 은행잎은 노랑고): 1.00000000000000000000000000000000  
P(은행잎은 노랑고): 0.00814814814814741827930788531  
P(노랑고): 0.00025460992907801416286137086153  
weighted probability of 빨강고 은행잎은 노랑고: 0.70165509062253739180903266969835

>> TRIGRAM INTERPOLATION:  
PROBABILITY: 1.3800156119401195e-13  
PERPLEXITY: 139.10817939561127

## 1.3 결과 확인

- (0.5, 0.3, 0.2)
  - 확률: 7.096168276562839e-14
  - 복잡도: 155.41610326068232
- (0.7, 0.2, 0.1)
  - 확률: 1.3800156119401195e-13
  - 복잡도: 139.10817939561127
- 세 모델 중 trigram 기준 확률이 항상 가장 높기 때문에 trigram 에 더 큰 가중치를 줄수록 확률 또한 높아지는 것을 확인할 수 있다.