
Homework1

Prepared By:
Humaira Qadeer

PART I:

We are given the following corpus, modified from the one in the chapter:

<s> I am Sam </s>

<s> Sam I am </s>

<s> I am Sam </s>

<s> I do not like green eggs and Sam </s>

Using a bigram language model with add-one smoothing, what is $P(\text{Sam} \mid \text{am})$? Include *<s>* and *</s>* in your counts just like any other token

$\text{Count}(\text{am Sam}) = 2$

$\text{Count}(\text{am}) = 3$

$V = 11$

$$P(\text{Sam} \mid \text{am}) = \frac{\text{Count}(\text{am Sam}) + 1}{\text{Count}(\text{am}) + V} = \frac{2+1}{3+11} = \frac{3}{14}$$

PART II:

1. The total unique word types in training corpus including the end of sentence padding symbol *</s>* and the unknown token *<unk>*, and disregarding the start of sentence padding symbol *<s>* is equal to 39502.

2. Disregarding the start of sentence padding symbol <s>, the total number of tokens in the training corpus is equal to 221290.
3. Before mapping the unknown words to <unk> in training and test data,
 - a the percentage of word tokens in the test corpus that did not occur in the training corpus is 1.873%
 - b the percentage of word types in the test corpus that did not occur in the training corpus is 3.929%
4. Before mapping the unknown words to <unk> in training and test data,
 - a the percentage of bigram tokens in the test corpus that did not occur in the training corpus is 25.4023%
 - b the percentage of bigram tokens in the test corpus that did not occur in the training corpus is 28.194%

5.

Unigram Model

- $\text{Log}_2(P(I)) = -8.400191146154114$
- $\text{Log}_2(P(\text{look})) = -11.982303986215943$
- $\text{Log}_2(P(\text{forward})) = -12.375607193140258$
- $\text{Log}_2(P(\text{to})) = -5.540690135458277$
- $\text{Log}_2(P(\text{hearing})) = -13.505537497185392$
- $\text{Log}_2(P(\text{your})) = -10.953683308276174$
- $\text{Log}_2(P(\text{reply})) = -17.623534706583843$
- $\text{Log}_2(P(.)) = -4.811868102497583$
- $\text{Log}_2(P(</s>)) = -4.62532894422938$
- $\text{Log}_2(P(S)) = -89.81874501974097$

Bigram Model

- $\text{Log}_2(P(<s>|i)) = -5.631088893700847$
- $\text{Log}_2(P(i|\text{look})) = -8.875420257009424$
- $\text{Log}_2(P(\text{look}|\text{forward})) = -4.345774836841731$
- $\text{Log}_2(P(\text{forward}|\text{to})) = -2.2644156362321546$
- $\text{Log}_2(P(\text{to}|\text{hearing})) = -13.2203480948755$
- $\text{Log}_2(P(\text{hearing}|\text{your})) = \text{undefined}$
- $\text{Log}_2(P(\text{your}|\text{reply})) = \text{undefined}$
- $\text{Log}_2(P(.|</s>)) = \text{undefined}$
- $\text{Log}_2(P(S)) = \text{Undefined}$

The log_2 probability for the sentence using the bigram model was undefined because the probability of $P(\text{hearing}|\text{your})$, $P(\text{your}|\text{reply})$, $P(.|</s>)$ equated to 0 which led to the log probability being undefined.

Bigram Model Smoothed

- $\text{Log}_2 (P(<s> | i)) = -6.155272148946965$
- $\text{Log}_2 (P(i | \text{look})) = -11.584868572205163$
- $\text{Log}_2 (P(\text{look} | \text{forward})) = -10.482195722024208$
- $\text{Log}_2 (P(\text{forward} | \text{to})) = -8.825392481587325$
- $\text{Log}_2 (P(\text{to} | \text{hearing})) = -13.827425410027113$
- $\text{Log}_2 (P(\text{hearing} | \text{your})) = -15.276596985010402$
- $\text{Log}_2 (P(\text{your} | \text{reply})) = -15.309973651164077$
- $\text{Log}_2 (P(\text{reply} | .)) = -15.270039765748985$
- $\text{Log}_2 (P(. | </s>)) = -0.6693438414242925$
- $\text{Log}_2 P(S) = -97.40110857813853$

6. Compute the perplexity of the sentence above under each of the models.

- Perplexity under Unigram Model : 1009.8046830192555
- Perplexity under Bigram Model: undefined
- Perplexity under Bigram Smoothed Model: 1810.7170379443896

The perplexity of the bigram model was undefined due to the probability of the sentence being undefined however the perplexity under the unigram model was less than the perplexity of the Bigram model with smoothing.

The unigram model assigned a higher probability to the occurrence of words which therefore resulted in the perplexity under the model to be lower where in contrast the smoothed bigram model assigned an increased probability to unseen events. In relation to the complexity of the models, the unigram model had high bias and low variance whereas the bigram smoothed model had a higher variance due to the probability mass of the zero probabilities allotting more probability space.

Compute the perplexity of the entire test corpus under each of the models. Discuss the differences in the results you obtained.

- Perplexity under Unigram Model : 1079.2881373075008
- Perplexity under Bigram Model: undefined
- Perplexity under Bigram Smoothed Model: 2640.1795213444884

The perplexity of the bigram model was undefined due to the probability of the sentence being undefined however the perplexity under the unigram model was less than the perplexity of the Bigram model with smoothing. The perplexity under the bigram smoothed model was much higher because of the weight it added to unseen events, redistributing the probabilities, and causing them to flatten. The unigram model is considered to be a better model based on its low perplexity in comparison to the other models.