

For 1.3 Questions.

It's hard to debug the program with very big training data.
Make a small training file, 3-4 sentences, a few words each sentence.

Do it by hand
To check logic error.

Part I - 3.4

$$P(\text{Sam} | \text{am}) = c(\text{am}, \text{Sam}) / c(\text{am})$$

$$\begin{array}{lcl} \text{am Sam} & 2 & = 2 / 3 \\ \text{am *} & 3 & \Downarrow \text{Add one smoothing} \\ V = 10 & & = (2+1) / (3+10) \\ & & = 3/13 \\ & & = 0.23077 \end{array}$$

</s>

I

am

Sam

do

not

like

green

eggs

and

1.3 Questions.

1. 41738

2. 2468210

3. Word token: 1.6612 %
Word type: 3.6058 %

4. Word token: 22.3679 %
Word type: 99.9763 %

5. S = i look forward to hearing your reply. </s>

Unigram params: i look forward to hearing your reply. </s>

$$P(S) = \log_2 P(i) + \log_2 P(\text{look}) + \log_2 P(\text{forward}) + \log_2 P(\text{to}) \\ + \log_2 P(\text{hearing}) + \log_2 P(\text{your}) + \log_2 P(\text{reply}) + \log_2 P(\cdot) \\ + \log_2 P(\text{</s>})$$

$$= -8.4510 - 12.0326 - 12.4036 - 5.3973 - 13.5850 \\ - 11.0432 - 17.5919 - 4.8689 - 4.6827 \\ = -90.2561$$

bigram params: i look, look forward, forward to, to hearing, hearing your, your reply, reply. , </s>

$$\log_2 \frac{\text{Count}(i, \text{look})}{\text{Count}(i)}$$

$$P(S) = \log_2 P(i \text{ look}) + \log_2 P(\text{look forward}) + \log_2 P(\text{forward to}) \\ + \log_2 P(\text{to hearing}) + \log_2 P(\text{hearing your}) + \log_2 P(\text{your reply}) \\ + \log_2 P(\text{reply} \cdot) + \log_2 P(\cdot \text{</s>})$$

p(parameter)
of these three
are 0.

The result of = -28.5463

the sentence under

bigram is undefined. But if we treat $\log_2 P(\cdot)$ of these taken as 0, we get ...

bigram with add one smoothing

parameters: i look, look forward, forward to, to hearing,
hearing your, your reply, reply ., . </s>

$$P(s) = \log_2 P(i \text{ look}) + \log_2 P(\text{look forward}) + \log_2 P(\text{forward to}) \\ + \log_2 P(\text{to hearing}) + \log_2 P(\text{hearing your}) + \log_2 P(\text{your reply}) \\ + \log_2 P(\text{reply} \cdot) + \log_2 P(\cdot </s>)$$

$$= \log_2 \frac{\text{count}(i, \text{look}) + 1}{\text{count}(i) + |V|^2} + \dots$$

$$= -26.6982 - 25.5689 - 24.0400 - 27.8909 \\ - 30.6982 - 30.6982 - 30.6982 - 14.3594 \\ = -210.6521$$

6. $s = i \text{ look forward to hearing your reply} \cdot </s>$

$$\text{Perplexity} = 2^{-L} \quad L = \frac{1}{M} \sum_{i=1}^M \log P(s_i)$$

$$M = 9$$

$$L = \frac{1}{9} \cdot (\log \text{prob})$$

$$\text{perplexity} = 2^{-\frac{1}{9} \cdot \log \text{prob of each model}}$$

$$\text{unigram ppl} = 2^{-\frac{1}{9} \cdot (-89.7404)} = 1003.7300$$

$$\text{bigram ppl} = 2^{-\frac{1}{9} \cdot (-28.5463)} = 9.0118$$

$$\text{bigram with add one smoothing ppl} = 2^{-\frac{1}{9} \cdot (-210.6521)} = 11113338.12$$

$\langle s \rangle$ included since Q7 not mentioned.

7. log prob

unigram $\bar{z} = -28592.7144$

bigram $\bar{z} = -13763.2607$

bigram w/ add one $\bar{z} = -72145.7616$

Perplexity / unigram $\bar{z} = 1283.6244$

bigram $\bar{z} = 31.3519$

bigram w/ add one $\bar{z} = 69707345.6550$

Inputs:

<s> A unigram maximum likelihood model . </s>

<s> A bigram maximum likelihood model . </s>

<s> A bigram model with Add-One smoothing . </s>

V = { a 3
unigram 1
maximum 2
likelihood 2
model 3
bigram 2
</s> 3
• 3
with 1
add-one 1
smoothing 1
}

$$\text{Total} = 7 + 7 + 8 = 22$$

$$\text{unigram: } p(a) = 3/22$$

$$p(\text{unigram}) = 1/22$$

$$p(\text{maximum}) = 2/22$$

$$p(\text{likelihood}) = 2/22$$

$$p(\text{model}) = 3/22$$

$$p(\text{bigram}) = 2/22$$

$$p(\text{•}) = 3/22$$

$$p(\text{with}) = 1/22$$

$$p(\text{add-one}) = 1/22$$

$$p(\text{smoothing}) = 1/22$$

bigram: p(a