

CS 124/LING 180/LING 280
From Languages to Information
Week 2: Group Exercises on Language Modeling
Winter 2019

Dan Jurafsky

Wednesday, January 23, 2019

Part 2: Group Exercise

We are interested in building a language model over a language with three words: A, B, C. Our training corpus is

AAABACBABBCCACBCC

1. First train a unigram language model using maximum likelihood estimation. What are the probabilities?
Reminder: We don't need start or end tokens for training a unigram model, since the context of each word doesn't matter. Thus, we won't add any special tokens to our corpus for now.

Answer:

$$P(A) = \frac{6}{18}$$

$$P(B) = \frac{6}{18}$$

$$P(C) = \frac{6}{18}$$

2. Next train a bigram language model using maximum likelihood estimation. For this problem, we'll add both a start token, $\langle start \rangle$, and an end token, $\langle end \rangle$, at the end of the string. These allow us to model the probability of the sentence starting with a particular letter and ending after a particular letter. Fill in the probabilities below. Leave your answers in the form of a fraction.

Answer:

Our corpus now becomes: $\langle start \rangle AAABACBABBCCACBCC \langle end \rangle$

$$\begin{aligned}
P(A | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle A)}{\text{count}(\langle start \rangle)} = \frac{1}{1} \\
P(A|A) &= \frac{\text{count}(AA)}{\text{count}(A)} = \frac{2}{6} \\
P(A|B) &= \frac{\text{count}(BA)}{\text{count}(B)} = \frac{2}{6} \\
P(A|C) &= \frac{\text{count}(CA)}{\text{count}(C)} = \frac{1}{6} \\
P(A | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle A)}{\text{count}(\langle end \rangle)} = \frac{0}{1} \\
P(B | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle B)}{\text{count}(\langle start \rangle)} = \frac{0}{1} \\
P(B|A) &= \frac{\text{count}(AB)}{\text{count}(A)} = \frac{2}{6} \\
P(B|B) &= \frac{\text{count}(BB)}{\text{count}(B)} = \frac{2}{6} \\
P(B|C) &= \frac{\text{count}(CB)}{\text{count}(C)} = \frac{2}{6} \\
P(B | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle B)}{\text{count}(\langle end \rangle)} = \frac{0}{1} \\
P(C | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle C)}{\text{count}(\langle start \rangle)} = \frac{0}{1} \\
P(C|A) &= \frac{\text{count}(AC)}{\text{count}(A)} = \frac{2}{6} \\
P(C|B) &= \frac{\text{count}(BC)}{\text{count}(B)} = \frac{2}{6} \\
P(C|C) &= \frac{\text{count}(CC)}{\text{count}(C)} = \frac{2}{6} \\
P(C | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle C)}{\text{count}(\langle end \rangle)} = \frac{0}{1} \\
P(\langle end \rangle | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle \langle end \rangle)}{\text{count}(\langle start \rangle)} = \frac{0}{1} \\
P(\langle end \rangle | A) &= \frac{\text{count}(A \langle end \rangle)}{\text{count}(A)} = \frac{0}{6} \\
P(\langle end \rangle | B) &= \frac{\text{count}(B \langle end \rangle)}{\text{count}(B)} = \frac{0}{6} \\
P(\langle end \rangle | C) &= \frac{\text{count}(C \langle end \rangle)}{\text{count}(C)} = \frac{1}{6} \\
P(\langle end \rangle | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle \langle end \rangle)}{\text{count}(\langle end \rangle)} = \frac{0}{1}
\end{aligned}$$

3. Now evaluate your language models on the test corpus:

ABACABB

What is the perplexity of the unigram language model evaluated on this corpus? Since we didn't add any special start/end tokens when we were training our unigram language model, we won't add any when we evaluate the perplexity of the unigram language model, either, so that we're consistent.

Answer:

$$\sqrt[7]{\left(\frac{1}{3}\right)^7} = \left(\frac{1}{3}\right)^{-1} = 3 \quad (1)$$

What is the perplexity of the bigram language model evaluated on this corpus? Since we added a start and end token when we were training our bigram model, we'll add them to this corpus again before we evaluate perplexity.

Answer: In this case, the corpus becomes: <start>ABACABB<end>

The perplexity is therefore:

$$\sqrt[8]{P(A|<start>)P(B|A)P(A|B)P(C|A)P(A|C)P(B|A)P(B|B)P(<end>|B)} \quad (2)$$

Note that in our unsmoothed bigram model, $P(<end>|B)$ is 0, since <end> never occurs after B in our training set. Therefore, the probability of this corpus is 0, so its perplexity is infinite.

4. Now repeat everything above for add-1 smoothing.

Answer: Add-one smoothing is simple for the unigram model. We simply add 1 to each numerator, and 3 (the vocabulary size) to each denominator. For this particular corpus, the actual probabilities are the same as without smoothing; all unigram probabilities are equal to one-third. Note that this is not true for most corpuses, as smoothing usually does in fact change the probabilities in your language model.

$$\begin{aligned} P(A) &= \frac{7}{21} \\ P(B) &= \frac{7}{21} \\ P(C) &= \frac{7}{21} \end{aligned}$$

Add-one smoothing is slightly more complicated for the bigram model. Since we use both an end-token <end>, our vocabulary size is 4 in this case, so we add 1 to each numerator and 4 to each denominator. Note that the <start> token is not considered part of our vocabulary as we don't expect it to every appear again after the first token in the corpus.

$$\begin{aligned}
P(A | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle A) + 1}{\text{count}(\langle start \rangle) + 4} = \frac{2}{5} \\
P(B | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle B) + 1}{\text{count}(\langle start \rangle) + 4} = \frac{1}{5} \\
P(C | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle C) + 1}{\text{count}(\langle start \rangle) + 4} = \frac{1}{5} \\
P(\langle end \rangle | \langle start \rangle) &= \frac{\text{count}(\langle start \rangle \langle end \rangle) + 1}{\text{count}(\langle start \rangle) + 4} = \frac{1}{5} \\
P(A|A) &= \frac{\text{count}(AA) + 1}{\text{count}(A) + 4} = \frac{3}{10} \\
P(A|B) &= \frac{\text{count}(BA) + 1}{\text{count}(B) + 4} = \frac{3}{10} \\
P(A|C) &= \frac{\text{count}(CA) + 1}{\text{count}(C) + 4} = \frac{2}{10} \\
P(A | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle A) + 1}{\text{count}(\langle end \rangle) + 4} = \frac{1}{5} \\
P(B|A) &= \frac{\text{count}(AB) + 1}{\text{count}(A) + 4} = \frac{3}{10} \\
P(B|B) &= \frac{\text{count}(BB) + 1}{\text{count}(B) + 4} = \frac{3}{10} \\
P(B|C) &= \frac{\text{count}(CB) + 1}{\text{count}(C) + 4} = \frac{3}{10} \\
P(B | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle B) + 1}{\text{count}(\langle end \rangle) + 4} = \frac{1}{5} \\
P(C|A) &= \frac{\text{count}(AC) + 1}{\text{count}(A) + 4} = \frac{3}{10} \\
P(C|B) &= \frac{\text{count}(BC) + 1}{\text{count}(B) + 4} = \frac{3}{10} \\
P(C|C) &= \frac{\text{count}(CC) + 1}{\text{count}(C) + 4} = \frac{3}{10} \\
P(C | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle C) + 1}{\text{count}(\langle end \rangle) + 4} = \frac{1}{5} \\
P(\langle end \rangle | A) &= \frac{\text{count}(A \langle end \rangle) + 1}{\text{count}(A) + 4} = \frac{1}{10} \\
P(\langle end \rangle | B) &= \frac{\text{count}(B \langle end \rangle) + 1}{\text{count}(B) + 4} = \frac{1}{10} \\
P(\langle end \rangle | C) &= \frac{\text{count}(C \langle end \rangle) + 1}{\text{count}(C) + 4} = \frac{2}{10} \\
P(\langle end \rangle | \langle end \rangle) &= \frac{\text{count}(\langle end \rangle \langle end \rangle) + 1}{\text{count}(\langle end \rangle) + 4} = \frac{1}{5}
\end{aligned}$$

The unigram perplexity ends up being the same as before (since the

smoothed unigram probabilities are the same as the unsmoothed unigram probabilities for this specific corpus):

$$\sqrt[7]{\left(\frac{1}{3}\right)^7} = \left(\frac{1}{3}\right)^{-1} = 3 \quad (3)$$

Bigram perplexity is a little bit different with smoothing. We use the same formula as before, but this time our probability is no longer 0, so we can compute a finite perplexity:

$$\begin{aligned} & \sqrt[8]{P(A|<start>)P(B|A)P(A|B)P(C|A)P(A|C)P(B|A)P(B|B)P(<end>|B)} \\ &= \sqrt[8]{\frac{2}{5}\left(\frac{3}{10}\right)^5\frac{2}{10}\frac{1}{10}} \\ &\approx 3.88 \end{aligned}$$

5. What is the difference between using an UNK token (for unknown words) and smoothing? In what situations would you use one versus the other?

Answer: UNK tokens are used when unknown words are observed in the dataset. This can happen in both train and test datasets - for some language model applications, it is desirable to have a fixed vocabulary with the UNK token representing the presence of other words during training and testing. If that is the case, at training and test time, any token not belonging to the fixed vocabulary would be converted into an UNK token before counting even starts.

Smoothing, on the other hand, is a way to redistribute the probabilities of observed occurrences from the training dataset. For almost all smoothing methods, the events (unigrams, bigrams ... etc) with higher probability would be discounted and rare events would be boosted. This is a good remedy for small datasets, where the probability distribution between unigrams and bigrams may be skewed due to the small size. However, in some cases this also causes the higher probability events to be discounted too much. Smoothing could also potentially be used to deal with unknown words - essentially treating unknown words as tokens that appear 0 times in the training corpus with their probabilities estimated from the smoothing addition.