

CS 124/LING 180/LING 280  
From Languages to Information  
Week 2: Group Exercises on Language Modeling  
Winter 2017

Dan Jurafsky

Tuesday, January 24, 2017

## 1 Part 1: 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. So, we won't add any special tokens to our corpus for this part of the problem.

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 an end token,  $\langle end \rangle$ , at the end of the string, so that we can model the probability of the sentence ending after a particular word. If you chose to add a start token as well, that's fine too, but these solutions assume no start token. Fill in the probabilities below. Leave your answers in the form of a fraction.

Answer:

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

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

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 an end token when we were training our bigram model, we'll add an end token to this corpus again before we evaluate perplexity.

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

The perplexity is therefore:

$$\sqrt[7]{P(B|A)P(A|B)P(C|A)P(A|C)P(B|A)P(B|B)P(<end>|B)} \quad (2)$$

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 used an end-token,  $<end>$ , our vocabulary size is 4 in this case, so we add 1 to each numerator and 4 to each denominator.

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

5678uj65442322409

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):

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

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:

$$\begin{aligned}
& \sqrt[7]{P(B|A)P(A|B)P(C|A)P(A|C)P(B|A)P(B|B)P(\langle end \rangle | B)} \\
&= \sqrt[7]{\left(\frac{3}{10}\right)^5 \frac{2}{10} \frac{1}{10}} \\
&= \sqrt[7]{\frac{486}{10^7}} \\
&= \frac{10}{486^{\frac{1}{7}}} \\
&\approx 4.132
\end{aligned}$$

## 2 Problem 2: Challenge Problems

1. What is the difference between using an UNK token (for unknown words) and smoothing? 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, and use the UNK token to treat any other words both when 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 probability distribution between unigrams and bigrams may be skewed due to the small size. However, in some cases this also cause the higher probability events to be discounted too much. Smoothing could potentially be used to deal with unknown words - essentially treating unknown words as tokens that appear 0 times in the training corpus, and their probabilities are estimated from the smoothing method.

They are really disjoint methods that should be used concurrently in a complete language model. The language models that you build only deal with UNK tokens at test time, where you should think about how to check for UNK tokens, and what probability to assign to the UNK tokens. At train time however, make sure you already perform the desired smoothing to the probabilities, so that you don't have to be worrying about that at score time. Even though the score function follows directly under the train function in the language model code, you should really be thinking about them as separated processes, where you apply smoothing in train and check for UNK in test.

2. Suppose you build an interpolated trigram language model, with three weights  $\lambda_1$  for unigrams,  $\lambda_2$  for bigrams, and  $\lambda_3$  for trigrams. Normally we set these lambdas on a held-out set. Suppose

instead we set them on the training data. This will cause the lambdas to take on very unusual values. What will these lambdas look like? Why?

$$(\lambda_1, \lambda_2, \lambda_3) = (0, 0, 1) \quad (4)$$

**Answer:** An n-gram model with a larger value of n is always a better fit for the training data, so the lambda value corresponding to the trigram model will be set to 1, and the bigram and unigram models will have lambda values equal to 0 (in other words, the bigram and unigram models will be completely ignored by the interpolated model).

This isn't usually desirable behavior, because even though our trigram model might fit our training data very well, it might overfit, and might not generalize well to an unseen test set. This is why we use a held-out set, which wasn't seen in training, for setting the lambda values.

3. Show that if we estimate two bigram language models using unsmoothed relative frequencies (MLE), one from a text corpus and the second from the same corpus in reverse order, the models will assign the same probability to new sentences (when applied in forward and backward order respectively). Hint: First write out the entire equation for sentence probabilities in terms of counts.

**Answer:** Let the sentence be  $(w_1, w_2, w_3, \dots, w_n)$ . Begin and end are the start and end tokens. T is the total number of sentences.

For forward order, the probability of the sentence is:

$$\begin{aligned} & Pr(\text{Begin}, w_1) \prod_{i=1}^{n-1} Pr(w_{i+1} | \text{Prevword is } w_i) Pr(\text{End} | w_n) \\ &= \frac{\text{count}(\text{Begin}, w_1)}{T} \prod_{i=1}^{n-1} \frac{\text{count}(w_i, w_{i+1})}{\text{count}(w_i)} \frac{\text{count}(w_n, \text{End})}{\text{count}(w_n)} \end{aligned}$$

For backward order, the probability of the sentence is:

$$\begin{aligned} & Pr(w_n, \text{End}) \prod_{i=1}^{n-1} Pr(w_i | \text{Nextword is } w_{i+1}) Pr(\text{Begin} | w_i) \\ &= \frac{\text{count}(w_n, \text{End})}{T} \prod_{i=1}^{n-1} \frac{\text{count}(w_i, w_{i+1})}{\text{count}(w_{i+1})} \frac{\text{count}(\text{Begin}, w_i)}{\text{count}(w_i)} \end{aligned}$$

Both are equal because

$$\text{count}(w_1) \prod_{i=1}^{n-1} \text{count}(w_{i+1}) = \text{count}(w_n) \prod_{i=1}^{n-1} \text{count}(w_i) \quad (5)$$