

Department of Computer Science  
The University of Melbourne  
COMP90042 WEB SEARCH AND TEXT ANALYSIS (Semester 1, 2018)

Workshop exercises: Week 7

**Discussion**

1. For the following “corpus” of two documents:

1. how much wood would a wood chuck chuck if a wood chuck would chuck wood
2. a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood

- I’m going to show the frequencies of the ten different word uni-grams, as it will make life a little easier in a moment:

a	chuck	could	he	how	if	much	the	wood	would	Total
4	9	1	1	1	2	1	1	8	4	32

(a) Which of the following sentences: A: a wood could chuck; B: wood would a chuck; is more probable, according to:

i. An unsmoothed uni-gram language model?

- An unsmoothed uni-gram language model is simply based on the counts of words in the corpus. For example, out of the 32 tokens in the corpus, there were 4 instances of a, so  $P(a) = \frac{4}{32}$
- To find the probability of a sentence using this model, we simply multiply the probabilities of the individual tokens:

$$\begin{aligned} P(A) &= P(a)P(\text{wood})P(\text{could})P(\text{chuck}) \\ &= \frac{4}{32} \times \frac{8}{32} \times \frac{1}{32} \times \frac{9}{32} \approx 2.75 \times 10^{-4} \end{aligned}$$

$$\begin{aligned} P(B) &= P(\text{wood})P(\text{would})P(a)P(\text{chuck}) \\ &= \frac{8}{32} \times \frac{4}{32} \times \frac{4}{32} \times \frac{9}{32} \approx 1.10 \times 10^{-3} \end{aligned}$$

- Clearly sentence B has the greater likelihood, according to this model.

ii. A uni-gram language model, with Laplacian (“add-one”) smoothing?

- Recall that in add-one smoothing, for each probability, we add 1 to the numerator and the size of the vocabulary (in this case, 10) to the denominator. For example,  $P_L(a) = \frac{4+1}{32+10} = \frac{5}{42}$ .
- Everything else proceeds the same way:

$$\begin{aligned} P_L(A) &= P_L(a)P_L(\text{wood})P_L(\text{could})P_L(\text{chuck}) \\ &= \frac{5}{42} \times \frac{9}{42} \times \frac{2}{42} \times \frac{10}{42} \approx 2.89 \times 10^{-4} \end{aligned}$$

$$\begin{aligned} P_L(B) &= P_L(\text{wood})P_L(\text{would})P_L(a)P_L(\text{chuck}) \\ &= \frac{9}{42} \times \frac{5}{42} \times \frac{5}{42} \times \frac{10}{42} \approx 7.23 \times 10^{-4} \end{aligned}$$

- Notice that the probability of sentence A is larger using this model, because the probability of the unlikely `could` has increased. (The other probabilities have decreased). Sentence B is still more likely, however.
- iii. An unsmoothed bi-gram language model?
- This time, we're interested in the counts of pairs of word tokens. For example, the probability of the bi-gram `wood would` is based on the count of that sequence of tokens, divided by the count of `wood`:  $\frac{1}{8}$  (because only a single `wood` is followed by `would`).
  - We are also going to include sentence terminals, so that the first probability in sentence A is  $P(a|<s>) = \frac{1}{2}$  — because one of the two sentences in the corpus starts with a. We also need to predict  $P(</s>|chuck) = \frac{0}{9}$  — because none of the 9 `chucks` are followed by the end of the sentence.
  - Now, we can substitute:

$$\begin{aligned}
 P(A) &= P(a|<s>)P(\text{wood}|a)P(\text{could}|\text{wood})P(\text{chuck}|\text{could})P(</s>|\text{chuck}) \\
 &= \frac{1}{2} \times \frac{4}{4} \times \frac{0}{8} \times \frac{1}{1} \times \frac{0}{9} = 0 \\
 P(B) &= P(\text{wood}|<s>)P(\text{would}|\text{wood})P(a|\text{would})P(\text{chuck}|a)P(</s>|\text{chuck}) \\
 &= \frac{0}{2} \times \frac{1}{8} \times \frac{1}{4} \times \frac{0}{4} \times \frac{0}{9} = 0
 \end{aligned}$$

- Because there is a zero-probability element in both of these calculations, they can't be nicely compared, leading us to instead consider:
- iv. A bi-gram language model, with Laplacian smoothing?
- We do the same idea as uni-gram add-one smoothing, but now the vocabulary size increases by one (because we're also predicting `</s>`; we need to do this to ensure that the probabilities of all the events that can follow a given token still sum to 1).

$$\begin{aligned}
 P_L(A) &= P_L(a|<s>)P_L(\text{wood}|a)P_L(\text{could}|\text{wood})P_L(\text{chuck}|\text{could})P_L(</s>|\text{chuck}) \\
 &= \frac{2}{13} \times \frac{5}{15} \times \frac{1}{19} \times \frac{2}{12} \times \frac{1}{20} \approx 2.25 \times 10^{-5} \\
 P_L(B) &= P_L(\text{wood}|<s>)P_L(\text{would}|\text{wood})P_L(a|\text{would})P_L(\text{chuck}|a)P_L(</s>|\text{chuck}) \\
 &= \frac{1}{13} \times \frac{2}{19} \times \frac{2}{15} \times \frac{1}{15} \times \frac{1}{20} \approx 3.60 \times 10^{-6}
 \end{aligned}$$

- This time, sentence A has the greater likelihood, mostly because of the common bi-gram `a wood`.
- v. An unsmoothed tri-gram language model?
- Same idea, longer contexts. Note that we now need two sentence terminals.

$$\begin{aligned}
 P(A) &= P(a|<s2> <s1>)P(\text{wood}|<s1> a) \cdots P(</s2>|\text{chuck} </s1>) \\
 &= \frac{1}{2} \times \frac{1}{1} \times \frac{0}{4} \times \frac{0}{0} \times \frac{0}{1} \times \frac{0}{0} = ? \\
 P(B) &= P(\text{wood}|<s2> <s1>)P(\text{would}|<s1> \text{wood}) \cdots P(</s2>|\text{chuck} </s1>) \\
 &= \frac{0}{2} \times \frac{0}{0} \times \frac{1}{1} \times \frac{0}{1} \times \frac{0}{0} \times \frac{0}{0} = ?
 \end{aligned}$$

- Given that the unsmoothed bi-gram probabilities were zero, that also means that the unsmoothed tri-gram probabilities will be zero. (Exercise for the reader: why?)
  - In this case, they aren't even well-defined, because of the  $\frac{0}{0}$  terms, but we wouldn't be able to meaningfully compare these numbers in any case.
- vi. A tri-gram language model, with Laplacian smoothing?
- The vocabulary size is now 12 (due to the two sentence terminals); everything else proceeds the same way:

$$\begin{aligned}
 P_L(A) &= P_L(a|<s2> <s1>)P_L(\text{wood}|<s1> a)\cdots P_L(</s2>|\text{chuck } </s1>) \\
 &= \frac{2}{14} \times \frac{2}{13} \times \frac{1}{16} \times \frac{1}{12} \times \frac{1}{13} \times \frac{1}{12} \approx 7.34 \times 10^{-7} \\
 P_L(B) &= P_L(\text{wood}|<s2> <s1>)P_L(\text{would}|<s1> \text{wood})\cdots P_L(</s2>|\text{chuck } </s1>) \\
 &= \frac{1}{14} \times \frac{1}{12} \times \frac{2}{13} \times \frac{1}{13} \times \frac{1}{12} \times \frac{1}{12} = 4.89 \times 10^{-7}
 \end{aligned}$$

- Notice that the problem of unseen contexts is now solved (they are just  $\frac{1}{12}$ ).
  - Sentence A has a slightly greater likelihood here, mostly because of the *a* at the start of one of the sentences (note that this will continue to be “seen” even at higher orders of *n*). You can also see that the numbers are getting very small, which is a good motivation for summing log probabilities (assuming no zeroes) rather than multiplying.
2. What does **back-off** mean, in the context of smoothing a language model? What does **interpolation** refer to?
- Back-off is a different smoothing strategy, where we incorporate lower-order *n*-gram models (in particular, for unseen contexts). For example, if we have never seen some tri-gram from our sentence, we can instead consider the bi-gram probability (at some penalty, to maintain the probability of all of the events, given some context, summing to 1). If we haven't seen the bi-gram, we consider the uni-gram probability. If we've never seen the uni-gram (this token doesn't appear in the corpus at all), then we need a so-called “0-gram” probability, which is a default for unseen tokens.
  - Interpolation is a similar idea, but instead of only “falling back” to lower-order *n*-gram models for unseen events, we can instead consider every probability as a linear combination of all of the relevant *n*-gram models, where the weights are once more chosen to ensure that the probabilities of all events, given some context, sum to 1.
3. Name the key differences and similarities between *n*-gram language models versus the log-bilinear language model and feed-forward neural language model.
- *n*-gram language models and LBL and FFNNs all share the same setup of using a Markov chain.

- They factorise the probability of a sentence into the probability of each word given the  $n - 1$  previous words. The models differ in how this word-based probability model (classifier) is formulated.
  - n-gram models can be considered a "feature-based" model where every ngram is a feature, with a corresponding weight (thus the "weights" for a bigram models form a matrix, for a trigram model a 3d tensor etc.)
  - LBL and FFNNLM use an embedding and un-embedding step, to limit the model size and force generalisation (e.g., to locate synonymous words near in vector space), along with more complex functions to couple context to the next word. The type of function is different in LBL vs FFLM, i.e., linear, vs a neural network.
  - Another key difference between n-grams and LBL/FFLM is that n-grams work over highly sparse data (1-hot word vectors, sparse parameter matrices where more entries are 0), while LBL/FFLM work over dense representations
4. What does **recurrent** mean in the context of a recurrent neural network (RNN) language model? How does the approach differ from a feed-forward language model?
- Recurrent means that model is structured such that it can be repeated applied for each item in a sequence. The recurrence in a RNN is over the hidden states, such that each as new input is fed into the model, this is used to formulate a new hidden state – as a non-linear transformation of the last hidden state, and the new input. In this way the approach can (in theory) represent long-distance phenomena in the sentence, over variable distances.
  - A FFLM instead assumes a fixed sized context, which can be applied using a "sliding window" over a sequence. The hidden state is a function of the n-1 inputs. There is no reuse of computation from previous applications to the sequence.