

Natural Language Processing and Machine Translation

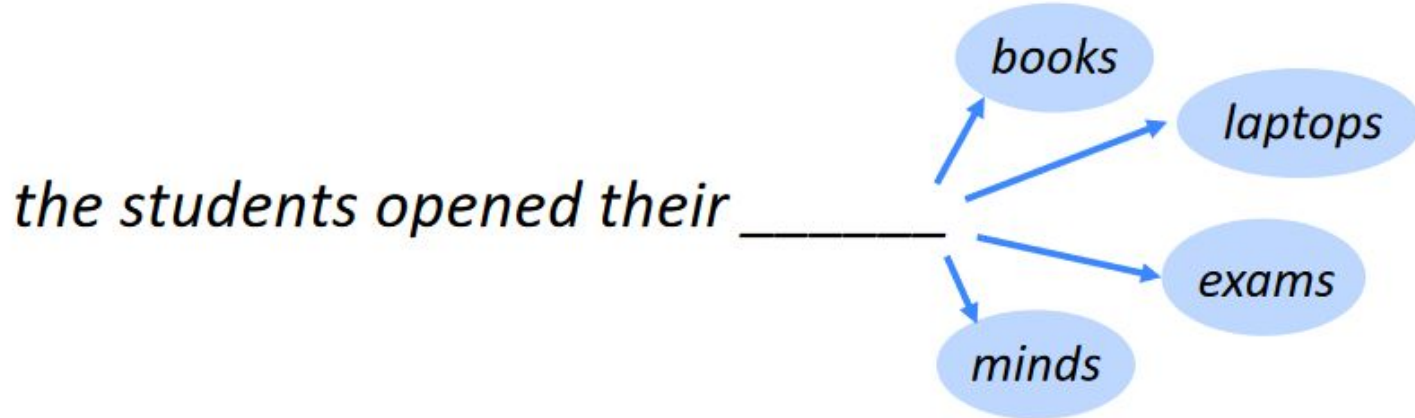
Language Models

Abhishek Koirala

M.Sc. in Informatics and
Intelligent Systems
Engineering

Introduction

- Use of various statistical and probabilistic techniques to determine the probability of a given sequence of words occurring in a sentence
- Analyze bodies of text data to provide a base for word predictions



<https://medium.com/@antonio.lopardo/the-basics-of-language-modeling-1c8832f21079>

N-gram

The cow jumps over the moon

Unigram/ 1-gram

The
cow
jumps
over
the
moon

Bigram/2-gram

The cow
cow jumps
jumps over
over the
the moon

3-gram

The cow jumps
cow jumps over
jumps over the
over the moon

4-gram

The cow jumps over
cow jumps over the
jumps over the moon

If X = Num of words in a given sentence K , the number of n -grams for sentence K would be:

$$Ngrams_K = X - (N - 1)$$

N-gram Language Models

Its water is so transparent that

$P(\text{the} | \text{its water is so transparent that})$.

One approach to calculate this using frequency approach

$$P(\text{the} | \text{its water is so transparent that}) = \frac{C(\text{its water is so transparent that the})}{C(\text{its water is so transparent that})}$$

Will this give us a good estimate in all possible scenarios ??

N-gram Language Models

Another way to do this is using **chain rule of probability**

$$p(w_1 \dots w_n) = p(w_1) \cdot p(w_2 | w_1) \cdot p(w_3 | w_1 w_2) \cdot p(w_4 | w_1 w_2 w_3) \dots p(w_n | w_1 \dots w_{n-1})$$

But this is again computationally expensive

We make this more simpler with an assumption:

- We approximate the context of the word w_k by looking at the last word of the context.
(**Markov Assumption**)

Eg. for bigram

$$p(w) = \prod_{i=1}^{k+1} p(w_i | w_{i-1})$$

N-gram language models

<s> I am a human </s>
<s> I am not a stone </s>
<s> I live in Lahore </s>

$$P(I|<S>) = C(<s>|I) / C(<s>) = 3/3 = 1$$

$$P(am|I) = C(I|am) / C(I) = 2/3$$

$$P(a|am) = C(am|a) / C(a) = 1/2$$

$$P(human|a) = C(a|human) / C(a) = 1/2$$

$$P(</s>|human) = C(human|</s>) / C(human) = 1$$

$$P(not|am) = C(am|not) / C(am) = 1/2$$

$$P(a|not) = C(not|a) / C(not) = 1$$

$$P(stone|a) = C(a|stone) / C(a) = 1/2$$

$$P(</s>|stone) = C(stone|</s>) / C(stone) = 1$$

$$P(live|I) = C(I|live) / C(I) = 1/3$$

$$P(in|live) = C(live|in) / C(live) = 1$$

$$P(Lahore|in) = C(in|Lahore) / C(in) = 1$$

$$P(</s>|Lahore) = C(Lahore|</s>) / C(Lahore) = 1$$

P(I am a human)

$$\begin{aligned} &= P(I|<s>) P(am|I) P(a|am) P(human|a) P(</s>|human) \\ &= 1 * 2/3 * 1/2 * 1/2 * 1 \\ &= 1/6 \end{aligned}$$

P(I am human)

$$\begin{aligned} &= P(I|<s>) P(am|I) P(human|am) P(</s>|human) \\ &= 1 * 2/3 * 0 * 1 \\ &= 0 \Rightarrow \text{Does this seem correct?} \end{aligned}$$

Laplace Smoothing

<s> I am a human </s>
<s> I am not a stone </s>
<s> I live in Lahore </s>

$$P(I|<S>) = C(<s>|I) / C(<s>) = 3/3 = 1$$

$$P(am|I) = C(I|am) / C(I) = 2/3$$

$$P(a|am) = C(am|a) / C(a) = 1/2$$

$$P(human|a) = C(a|human) / C(a) = 1/2$$

$$P(</s>|human) = C(human|</s>) / C(human) = 1$$

$$P(not|am) = C(am|not) / C(am) = 1/2$$

$$P(a|not) = C(not|a) / C(not) = 1$$

$$P(stone|a) = C(a|stone) / C(a) = 1/2$$

$$P(</s>|stone) = C(stone|</s>) / C(stone) = 1$$

$$P(live|I) = C(I|live) / C(I) = 1/3$$

$$P(in|live) = C(live|in) / C(live) = 1$$

$$P(Lahore|in) = C(in|Lahore) / C(in) = 1$$

$$P(</s>|Lahore) = C(Lahore|</s>) / C(Lahore) = 1$$

The solution to the problem of unseen N-grams is to re-distribute some of the probability mass from the observed frequencies to unseen N-grams. This is a general problem in probabilistic modeling called **smoothing**.

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$$

Using laplace smoothing (Vocab = 11)

P(I am human)

$$\begin{aligned} &= P(I|<s>) P(am|I) P(human|am) P(</s>|human) \\ &= (3+1)/(3+11) * (2+1)/(3+11) * (0+1)/(2+11) * (1+1)/(1+11) \\ &= 4/14 * 3/14 * 1/13 * 2/12 \\ &= 0.00078 \end{aligned}$$

Good Turing Discounting

- Re-estimate the amount of probability mass to assign N-gram with zero or low counts by looking at the number of N-grams with higher counts
- Use the count of things which are seen once to help estimates the count of things never seen.
- Let N_c be number of N-grams that occur c times
 - For bigrams, N_0 , is the number of bigrams of count 0, N_1 , is the number of bigrams with count 1, etc
- Revised count

$$c^* = (c + 1) \frac{N_{c+1}}{N_c}$$

