

## Statistical NLP 2019 - Assignment 1

Consider  $n$ -gram language modeling:

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | w_{i-n+1} \dots w_{i-1}) = P(w_i | w_{i-n+1}^{i-1})$$

let  $V$  denote the vocabulary size (typically around a million),  
 $n$  the  $n$ -gram order (typically around 5), and  $C$  the  
 number of tokens in the corpus used to compute the  
 relevant counts (typically around 1 billion).

1. Characterize the memory complexity of Kneser-Ney language models in big- $O$  notation. You should decide what variables are important to model. (Hint: think about how to store the relevant counts efficiently)

Solution: The key is to find a way to store the count of all  $(w', w)$  pairs for all  $k$ -grams for  $k = 1, 2, \dots, n$ , where  $w'$  is the history and  $w$  is the next word.

Consider storing the relevant counts in a matrix for a specific  $k$ -gram, the number of different histories (number of rows) is bounded by:

$\min(V^k, C)$ , because each token can take on  $V$  possible values, and the corpus length is  $C$ .

Since  $V$  is around a million, which is  $1 \times 10^6$ ,

$C$  is around a billion, which is  $1 \times 10^9$ , we have

$$\min(V^k, d(C)) = \begin{cases} V^k, & \text{if } k=1, \\ d(C), & \text{otherwise} \end{cases}$$

The number of different words (number of columns) is the vocabulary size, which is  $V$ .

Using a dense matrix representation, the space complexity of the matrix for a specific  $k$ -gram is  $O(VC)$  for  $k=2, \dots, n$ , and  $O(V^2)$  for  $k=1$ .

Thus the total space complexity is:

$$(n-1)O(VC) + O(V^2) = O(nVC)$$

However,  $O(VC)$  is around  $1 \times 10^{15}$ , so using the dense matrix representation may not work.

Instead, consider the sparse matrix representation, the space complexity of the matrix for a specific  $k$ -gram is the number of non-empty entries in the corresponding dense matrix, which is  $O(C)$  for  $k=1, 2, \dots, n$ .

Thus the total space complexity is:

$$nO(C) = O(nC),$$

which is around  $n \times 10^9$  for  $n$  around 5.

The sparse matrix representation might be feasible.



2. The term inference time refers to one call to the language model, i.e., computing  $P(W_i | W_{i-n+1}^{i-1})$  under the language model for one choice of  $W_{i-n+1}^{i-1}, W_i$ .

Characterize the inference time complexity of Kneser-Ney language models in big-O notation. (Hint: Think about which quantities can be cached efficiently to speed up inference time).

Solution: The recursive equation for computing  $P_{KN}$  is:

$$P_{KN}(W_i | W_{i-n+1}^{i-1}) = \frac{\max(C_{KN}(W_{i-n+1}^i) - d, 0)}{\sum_v C_{KN}(W_{i-n+1}^{i-1} v)} + \lambda(W_{i-n+1}^{i-1}) P_{KN}(W_i | W_{i-n+2}^{i-1})$$

Where  $C_{KN}(W') = \begin{cases} \text{count}(W'), & \text{for } W' = W_{i-n+1}^i, \\ \text{continuation count}(W'), & \text{otherwise} \end{cases}$

$\text{continuation count}(W') =$  the number of unique single word contexts for  $W'$ .

$$\lambda(W_{i-n+1}^{i-1}) = \frac{d}{\sum_v C_{KN}(W_{i-n+1}^{i-1} v)} |\{w: C(W_{i-n+1}^{i-1} w) > 0\}|$$

We can cache the following quantities in advance:

① For each  $(n, i)$  pair, store  $\sum_v C_{KN}(W_{i-n+1}^{i-1} v)$ .

This requires  $O(nc)$  space.

②  $|\{w: C(W_{i-n+1}^{i-1} w) > 0\}|$  for each  $(n, i)$  pair.

This requires  $O(nc)$  space.

We cannot store each  $(W_{i-n+1}^{i-1}, w^i)$  pair since this requires  $O(VC)$  space. Instead, to compute  $C_{KN}(W_{i-n+1}^i) = C_{KN}(W_{i-n+1}^{i-1} w^i)$ , we search in the row indexed by  $W_{i-n+1}^{i-1}$  for the column  $w^i$  in the corresponding  $n$ -gram count matrix. Indexing into an element in the sparse matrix takes time proportional to the logarithm of the length of its columns, which is  $O(\log V)$ , usually around  $\log_2(11 \times 10^4) \approx 20$ . This overhead is acceptable.

To sum up, if we cache  $\sum_v C_{KN}(W_{i-n+1}^{i-1} v)$  and  $|\{w: C(W_{i-n+1}^{i-1} w) > 0\}|$  for each  $(n, i)$  pair in advance, requiring  $O(nc)$  space overhead, we can compute  $\sum_v C_{KN}(W_{i-n+1}^{i-1} v)$  and  $\lambda(W_{i-n+1}^{i-1})$  in  $O(1)$  time. We can also compute  $\max(C_{KN}(W_{i-n+1}^i) - d, 0)$  in  $O(\log V)$  time.

Thus we have the recurrence:

$$T(n) = O(\log V) + T(n-1)$$

which gives  $T(n) = O(n \log V)$

So the time complexity for inference time is  $O(n \log V)$



3. Recall absolute discounting for  $n=2$ :

$$P_{\text{ad}}(w_i | w_{i-1}) = \frac{\max(c(w_{i-1}, w_i) - d, 0)}{\sum_w c(w_{i-1}, w)} + \alpha(w_{i-1}) \hat{p}(w_i)$$

where  $\hat{p}(w_i) = \frac{c(w_i)}{\sum_w c(w)}$  is the empirical unigram dist'n,

$d$  is the constant and  $\alpha$  is the left-over weight assigned to the lower order distribution:

$$\alpha(w_{i-1}) = \frac{d \times |\{w_i: c(w_{i-1}, w_i) > 0\}|}{\sum_w c(w_{i-1}, w)}$$

For  $n=2$  (a bigram model), give an example of a small vocabulary  $V$  and corpus  $C$  where absolute discounting does not preserve the marginal constraint, i.e.,

$$\hat{p}(w_i) \neq \sum_{w_{i-1}} P_{\text{ad}}(w_i | w_{i-1}) \cdot \hat{p}(w_{i-1}) \quad (1)$$

Solution:

Corpus:  $\langle s \rangle$  I am A I am B  $\langle /s \rangle$

Vocabulary:  $\{\langle s \rangle, \text{I}, \text{am}, \text{A}, \text{B}, \langle /s \rangle\}$

$$\hat{p}(\text{I}) = \frac{c(\text{I})}{\sum_w c(w)} = \frac{2}{8} = \frac{1}{4}$$

$$\begin{aligned} \text{let } \tilde{p}(\text{I}) &= \sum_{w_{i-1}} P_{\text{ad}}(\text{I} | w_{i-1}) \cdot \hat{p}(w_{i-1}) \\ &= P_{\text{ad}}(\text{I} | \langle s \rangle) \cdot \hat{p}(\langle s \rangle) + P_{\text{ad}}(\text{I} | \text{A}) \cdot \hat{p}(\text{A}) \\ &= \left( \frac{1-d}{1} + \frac{d \times 1}{1} \times \frac{1}{4} \right) \times \frac{1}{8} \times 2 \\ &= \frac{1}{4} - \frac{9}{16} d \end{aligned}$$

We have  $\frac{1}{4} = \hat{p}(\text{I}) \neq \tilde{p}(\text{I}) = \frac{1}{4} - \frac{9}{16} d$  for  $0 < d < 1$