

Slp3-chapter3-Exercises-20200428

3.1 Write out the equation for trigram probability estimation (modifying Eq. 3.11).

Now write out all the non-zero trigram probabilities for the I am Sam corpus

on page 33.

【1】

LaTeX equation :

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}w)}{C(w_{n-1})}$$

【2】

Because <s> I am Sam </s>, so :

$$\begin{aligned} P(<\text{s}> \text{ I am Sam } </\text{s}>) \\ = & P(\text{am} | <\text{s}> \text{ I }) P(\text{Sam} | \text{ I am }) P(</\text{s}> | \text{ am Sam }) \\ = & 0.5 \times 0.5 \times 1 \\ = & 0.25 \end{aligned}$$

3.2 Calculate the probability of the sentence i want chinese food. Give two probabilities, one using Fig. 3.2 and the ‘useful probabilities’ just below it on page 35, and another using the add-1 smoothed table in Fig. 3.6. Assume the additional add-1 smoothed probabilities $P(i|<\text{s}>) = 0.19$ and $P(</\text{s}>|\text{food}) = 0.40$.

【1】 2-grams

$$\begin{aligned} P(<\text{s}> \text{ i want chinese food } </\text{s}>) \\ = & P(\text{i} | <\text{s}>) P(\text{want} | \text{i}) P(\text{chinese} | \text{want}) P(\text{food} | \text{chinese}) P(</\text{s}> | \text{food}) \\ = & 0.25 \times 0.33 \times 0.0065 \times 0.52 \times 0.68 \\ = & 0.000189618 \end{aligned}$$

【2】 2-grams (add-1 smoothing)

$$\begin{aligned} P^*(<\text{s}> \text{ i want chinese food } </\text{s}>) \\ = & P^*(\text{i} | <\text{s}>) P^*(\text{want} | \text{i}) P^*(\text{chinese} | \text{want}) P^*(\text{food} | \text{chinese}) P^*(</\text{s}> | \text{food}) \\ = & 0.19 \times 0.21 \times 0.0029 \times 0.052 \times 0.40 \\ = & 0.00000240676 \end{aligned}$$

3.3 Which of the two probabilities you computed in the previous exercise is higher, unsmoothed or smoothed? Explain why.

Answer: Unsmoothed. The sharp change in counts and probabilities occurs because too much probability mass is moved to all the zeros.

3.4 We are given the following corpus, modified from the one in the chapter:

<s> I am Sam </s>

<s> Sam I am </s>

<s> I am Sam </s>

<s> I do not like green eggs and Sam </s>

Using a bigram language model with add-one smoothing, what is $P(\text{Sam} \mid \text{am})$? Include <s> and </s> in your counts just like any other token.

$$\begin{aligned} P_{\text{Laplace}}^*(\text{Sam} \mid \text{am}) &= \frac{C(\text{am Sam}) + 1}{C(\text{am}) + V} \\ &= \frac{2 + 1}{3 + 25} \\ &= \frac{3}{28} \\ &= 0.1071428 \end{aligned}$$

3.5 Suppose we didn't use the end-symbol </s>. Train an unsmoothed bigram grammar on the following training corpus without using the end-symbol </s>:

<s> a b

<s> b b

<s> b a

<s> a a

Demonstrate that your bigram model does not assign a single probability distribution across all sentence lengths by showing that the sum of the probability of the four possible 2 word sentences over the alphabet {a,b} is 1.0, and the sum of the probability of all possible 3 word sentences over the alphabet {a,b} is also 1.0.

3.6 Suppose we train a trigram language model with add-one smoothing on a given corpus. The corpus contains V word types. Express a formula for estimating $P(w_3|w_1, w_2)$, where w_3 is a word which follows the bigram (w_1, w_2) , in terms of various N -gram counts and V . Use the notation $c(w_1, w_2, w_3)$ to denote the number of times that trigram (w_1, w_2, w_3) occurs in the corpus, and so on for bigrams and unigrams.

22:22 4月28日周二 20%

备忘录 (trigram) $\textcircled{1}$

Formula : $P^*(w_3 | w_1, w_2)$

$$= \frac{C(w_1, w_2, w_3) + 1}{C(w_1, w_2) + V}$$

(bigram) $\textcircled{2}$

$$P^*(w_3 | w_2) = \frac{C(w_2, w_3) + 1}{C(w_2) + V}$$

Allgrams $\textcircled{3}$

$$P^*(w_n | w_{n-N+1}, w_{n-N+2}, \dots, w_{n-1})$$

$$= \frac{C(w_{n-N+1}, w_{n-N+2}, \dots, w_n) + 1}{C(w_{n-N+1}, w_{n-N+2}, \dots, w_{n-1}) + V}$$

3.7 We are given the following corpus, modified from the one in the chapter:

<s> I am Sam </s>

<s> Sam I am </s>

<s> I am Sam </s>

<s> I do not like green eggs and Sam </s>

If we use linear interpolation smoothing between a maximum-likelihood bigram model and a maximum-likelihood unigram model with $\lambda_1 = 0.5$

and $\lambda_2 = 0.5$

, what is $P(\text{Sam}|\text{am})$? Include <s> and </s> in your counts just like any other token.

$$\begin{aligned}\hat{P}(\text{sam}|\text{am}) &= \lambda_1 P(\text{sam}|\text{am}) + \lambda_2 P(\text{sam}) \\ &= 0.5 \times \frac{2}{3} + 0.5 \times 1 \\ &= \frac{5}{6}\end{aligned}$$

3.8 Write a program to compute unsmoothed unigrams and bigrams.

3.9 Run your n-gram program on two different small corpora of your choice (you might use email text or newsgroups). Now compare the statistics of the two corpora. What are the differences in the most common unigrams between the two? How about interesting differences in bigrams?

3.10 Add an option to your program to generate random sentences.

3.11 Add an option to your program to compute the perplexity of a test set.

3.12 Given a training set of 100 numbers consists of 91 zeros and 1 each of the other digits 1-9. Now we see the following test set: 0 0 0 0 0 3 0 0 0 0. What is the unigram perplexity?

23:47 4月28日周二

9%

备忘录



Fix 3.12 $\phi =$

a corpus contains 10 test set.

a test set has 10 words

$$\text{Perplexity}(\text{corpus}) = 2^{-\frac{1}{10} \sum_{i=1}^{10} \log_2 P(s_i)}$$

$$\text{由于 } P(s_i) = \prod_{j=1}^n P(w_j)$$

由于 0,0,0,0,0,0,3,0,0,0 中有 0,3.

$$\text{Fix } P(s_i) = P(3) \cdot P(0) \xrightarrow{\text{unigram-gram.}} \\ = \frac{1}{10} \cdot \frac{9}{10} = \frac{9}{100}$$

从上述推导，corpus 固定度为。

$$= 2^{-\frac{1}{10} \sum_{i=1}^{10} \log_2 \left(\frac{9}{100} \right)} \xrightarrow{\left(\frac{9}{10} \right)^2} (0.9)^2$$

$$= 2^{-\frac{1}{10} \cdot (20 + \sum_{i=1}^{10} (-1.7))}$$

$$= 2^{-\frac{1}{10} \times (20 - 17)}$$

$$= \sqrt[10]{2} \times 3$$

$$= 1.07 \times 3$$

