

Einführung in die Computerlinguistik

Hausaufgabe Language Models, Abgabe 18.05.2020

Laura Kallmeyer

SoSe 2020, Heinrich-Heine-Universität Düsseldorf

Aufgabe 1 Consider a toy example where the sentences are built over $\{a, b\}$ and our training data consists of the following sentences:

$\langle s \rangle a a b a a b b \langle /s \rangle$

$\langle s \rangle a b a b b a a b a a \langle /s \rangle$

$\langle s \rangle b b a a b a b a b a \langle /s \rangle$

$\langle s \rangle b b a b a b \langle /s \rangle$

1. Give the following bigram probabilities and their \log_2 values:

$$P(b|\langle s \rangle) \quad P(b|b) \quad P(\langle /s \rangle|b)$$

2. Using these values, the probability of a sequence $\langle s \rangle w_1 w_2 \dots w_n$ amounts to

$$P(\langle s \rangle w_1 w_2 \dots w_n) = 2^{\sum_{k=1}^n \log_2 P(w_k | w_{k-1})}$$

while the perplexity value of this sequence is

$$PP(\langle s \rangle w_1 w_2 \dots w_n) = \sqrt[n]{\frac{1}{2^{\sum_{k=1}^n \log_2 P(w_k | w_{k-1})}}} = \sqrt[n]{2^{\sum_{k=1}^n |\log_2 P(w_k | w_{k-1})|}}$$

Compute both, probability and perplexity, of the following sequences:

(1) $\langle s \rangle b \langle /s \rangle$

(2) $\langle s \rangle b b b \langle /s \rangle$

Solution:

1. $P(b|\langle s \rangle) = \frac{1}{2}$, $\log_2 \frac{1}{2} = -1$ $P(b|b) = \frac{4}{16}$, $\log_2 \frac{4}{16} = -2$ $P(\langle /s \rangle|b) = \frac{1}{8}$, $\log_2 \frac{1}{8} = -3$

2. (1): probability $2^{-1-3} = 2^{-4} = \frac{1}{2^4} = \frac{1}{16}$, perplexity $\sqrt[4]{2^4} = 2$

(2): probability $2^{-1-2-2-3} = 2^{-8} = \frac{1}{2^8} = \frac{1}{612}$, perplexity $\sqrt[4]{2^8} = 2^2 = 4$

Aufgabe 2 Assume that we want to obtain a bigram language model with Laplace Smoothing. Our training data are the following two sentences:

$\langle s \rangle b a a b a \langle /s \rangle$

$\langle s \rangle a b a b a b a b b \langle /s \rangle$

(The size of the vocabulary is 3 since we have three possible second elements of our bigrams, a , b , and $\langle /s \rangle$.)

This time, we compute using the probabilities, not their log values.

1. What are the following probabilities according to this model?

(a) $P(a|\langle s \rangle)$

(b) $P(b|a)$

(b) $P(a|a)$

(b) $P(\langle /s \rangle|b)$

2. Give probability and perplexity of the following input sentence:

$\langle s \rangle a a b \langle /s \rangle$

Solution:

1. (a) $P(a|\langle s \rangle) = \frac{2}{5}$ (b) $P(b|a) = \frac{6}{10} = \frac{3}{5}$ (c) $P(a|a) = \frac{2}{10} = \frac{1}{5}$
 (d) $P(\langle /s \rangle|b) = \frac{2}{10} = \frac{1}{5}$
2. Probability: $\frac{2}{5} \cdot \frac{1}{5} \cdot \frac{3}{5} \cdot \frac{1}{5} = \frac{6}{625}$
 Perplexity: $\sqrt[4]{\frac{625}{6}}$

Aufgabe 3 Consider again the same training data as in the preceding question.

1. Modify the training data adopting a treatment of unknown words that changes the first occurrence of every word into $\langle \text{UNK} \rangle$.

Give the following bigram probabilities estimated by this model:

$$P(a|\langle \text{UNK} \rangle) \quad P(\langle \text{UNK} \rangle|a) \quad P(\langle \text{UNK} \rangle|\langle s \rangle) \quad P(a|\langle s \rangle)$$

2. Calculate the probability of the two sequences

- (3) a. $\langle s \rangle a c$
 b. $\langle s \rangle c a$

according to this model.

3. Now add Laplace Smoothing. What are the new values of our four bigram probabilities and the probabilities of (3-a) and (3-b) in this new model?

Solution:

1. Modified training data:

$\langle s \rangle \langle \text{UNK} \rangle \langle \text{UNK} \rangle a b a \langle /s \rangle$
 $\langle s \rangle a b a b a b a b \langle /s \rangle$

$$P(a|\langle \text{UNK} \rangle) = \frac{1}{2} \quad P(\langle \text{UNK} \rangle|a) = 0 \quad P(\langle \text{UNK} \rangle|\langle s \rangle) = \frac{1}{2} \quad P(a|\langle s \rangle) = \frac{1}{2}$$

2. $P(\langle s \rangle ac) = \frac{1}{2} \cdot 0 = 0$ $P(\langle s \rangle ca) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$

3. With additional Laplace Smoothing:

$$P(a|\langle \text{UNK} \rangle) = \frac{2}{6} \quad P(\langle \text{UNK} \rangle|a) = \frac{1}{10} \quad P(\langle \text{UNK} \rangle|\langle s \rangle) = \frac{2}{6} \quad P(a|\langle s \rangle) = \frac{2}{6}$$

$$P(\langle s \rangle ac) = \frac{1}{3} \cdot \frac{1}{10} = \frac{1}{30} \quad P(\langle s \rangle ca) = \frac{1}{3} \cdot \frac{1}{3} = \frac{1}{9}$$